



# Minimizing a differentiable function over a differential manifold

D. Gabay

## ► To cite this version:

D. Gabay. Minimizing a differentiable function over a differential manifold. [Research Report] RR-0009, INRIA. 1980. inria-00076552

**HAL Id: inria-00076552**

**<https://inria.hal.science/inria-00076552>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Rapports de Recherche

N° 9

**MINIMIZING  
A DIFFERENTIABLE FUNCTION  
OVER  
A DIFFERENTIAL MANIFOLD**

**Daniel GABAY**

**Février 1980**

Institut National  
de Recherche  
en Informatique  
et en Automatique

Domaine de Voluceau  
Rocquencourt  
BP 105 78150 Le Chesnay  
France  
Tél. 954 90 20

# MINIMIZING A DIFFERENTIABLE FUNCTION

## OVER A DIFFERENTIAL MANIFOLD

by

Daniel GABAY

### ABSTRACT

We present in this paper numerical methods for optimization problems constrained by nonlinear equalities. In the first part we define intrinsically the gradient field of the objective function on the constraint manifold and analyze descent methods along geodesics, including the gradient projection and reduced gradient methods for special choices of coordinate systems. In particular we generalize the Quasi-Newton methods and establish their superlinear convergence.

In the second part we present an efficient approximation of the Quasi-Newton method along geodesics where the feasibility of successive iterates is improved instead of being enforced. The method uses an exact penalty function and converges globally and superlinearly. Its relation with multipliers methods and recursive quadratic programming methods is investigated, and its advantages evidenced.

### RESUME

On présente dans ce rapport des méthodes numériques pour les problèmes d'optimisation sous contraintes d'égalité non-linéaires. Dans la première partie on définit intrinsèquement le champ de gradient de la fonction objectif sur la variété définie par les contraintes et on analyse des méthodes de descente le long de géodésiques ; elles incluent les méthodes classiques du gradient projeté et du gradient réduit qui correspondent à des choix particuliers de systèmes de coordonnées. On généralise notamment les méthodes quasi-Newtoniennes dont on établit la convergence superlinéaire.

Dans la seconde partie on présente une approximation efficace de la méthode quasi-Newtonienne le long de géodésiques où l'admissibilité des itérés successifs est seulement améliorée au lieu d'être strictement imposée. La méthode utilise

une fonction de pénalisation exacte et converge globalement avec une vitesse superlinéaire. On étudie aussi les relations entre cette nouvelle méthode et les méthodes de multiplicateurs ainsi que les méthodes utilisant des programmes quadratiques rékursifs mettant ainsi ses avantages en évidence.

MINIMIZING A DIFFERENTIABLE FUNCTION  
OVER A DIFFERENTIAL MANIFOLD  
PART I : DESCENT METHODS ALONG GEODESICS  
AND PRACTICAL IMPLEMENTATION\*

Daniel GABAY  
Laboratoire d'Analyse Numérique, Université P. et M. Curie (PARIS VI)

and

Institut National de Recherche en Informatique et Automatique  
Domaine de Voluceau, 78150 Le Chesnay (France)

ABSTRACT

To generalize the descent methods of unconstrained optimization to the constrained case, we define intrinsically the gradient field of the objective function on the constraint manifold and analyze descent methods along geodesics, including the gradient projection and reduced gradient methods for special choices of coordinate systems. In particular we generalize the Quasi-Newton methods and establish their superlinear convergence ; we show that they only require the updating of a reduced size matrix. In practice the geodesic search is approximated by a tangent step followed by a constraints restoration or by a simple arc search again followed by a restoration step.

\* presented at the 10<sup>th</sup> International Symposium on Mathematical Programming, Montréal, August 1979.

# 1 - INTRODUCTION

This paper shares with a previous article by the author together with D. Luenberger (Ref. 1) the purpose of extending the well-known gradient-related methods for the unconstrained minimization of a real-valued function on  $\mathbb{R}^n$  to the nonlinearly constrained problem

$$(1.1) \quad \text{Min } \{f(x) \mid x \in \mathbb{R}^n : c_i(x) = 0, i = 1, 2, \dots, m\},$$

where  $m \leq n$ ,  $f$  and  $c_i$  are real-valued functions on  $\mathbb{R}^n$  and assumed to be  $\mathcal{C}^\sigma$  differentiable ( $\sigma \geq 2$  unless otherwise specified). Denoting by  $c$  the map from  $\mathbb{R}^n$  into  $\mathbb{R}^m$  of component  $c_i$ , we also assume that the following regularity assumption holds : 0 is a regular value of the map  $c$ , i.e. the Jacobian  $c'(x) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$  of  $c$  at  $x$  is of full rank  $m$  for all  $x \in C = c^{-1}(0)$ .

In (Ref. 1) this program was fulfilled in the following way. The regularity assumption implies that there exists an  $m \times m$  submatrix of the Jacobian matrix

$A_x$  of  $c$  at  $x \in C$  which is non-singular; say  $B = \left( \frac{\partial c_i}{\partial x_j} \right)$ ,  $i, j \in I = \{1, \dots, m\}$  and let  $J = N - I$ .

Identify  $\mathbb{R}^n$  with  $\mathbb{R}^m \times \mathbb{R}^{n-m}$  and set  $x = (x_I, x_J)$ ,  $A_x = [B, D]$ .

According to the implicit function theorem, there exist a neighborhood

$U_x = W \times V$  of  $x = (x_I, x_J)$  in  $\mathbb{R}^m \times \mathbb{R}^{n-m}$  and a  $\mathcal{C}^\sigma$  map  $\psi_x : V \rightarrow W$  such that

$y_I = \psi_x(y_J)$  if and only if  $c(y_I, y_J) = 0$ . The restriction of  $f$  to  $U_x \cap C$  can thus be represented by a  $\mathcal{C}^\sigma$  real-valued function  $\phi_x : V \subset \mathbb{R}^{n-m} \rightarrow \mathbb{R}$  defined by

$$(1.2) \quad \phi_x(z) = f(\psi_x(z), z).$$

Suppose now that  $x^* \in C$  is a local minimizer of  $f$  on  $C$  in the neighborhood  $U_x \cap C$  of  $x$ . Then  $x^* = (x_I^*, x_J^*)$  with  $x_I^* = \psi_x(x_J^*)$  and  $x_J^*$  is a local minimizer of the (unconstrained) minimization problem on the open set  $V$

$$(1.3) \quad \text{Min}_{z \in V} \phi_x(z).$$

Applying the gradient-related methods to the unconstrained reduced problem (1.3) we construct a sequence of approximations  $\{z^k\}$  converging to  $x_J^*$  defined iteratively by

$$(1.4) \quad z^{k+1} = z^k + \tau_k p^k,$$

where  $p^k$  is a descent direction based on  $g^k$  the gradient of  $\varphi_x$  at  $z^k$ , and  $t_k$  a stepsize selected by a line search of the function  $j(t) = \varphi_x(z^k + t p^k)$  for  $t > 0$ . In the original space  $\mathbb{R}^n$  iteration (1.4) corresponds to the scheme

$$(1.5 \text{ a}) \quad x_J^{k+1} = x_J^k + t_k p^k,$$

$$(1.5 \text{ b}) \quad x_I^{k+1} = \psi(x_J^{k+1}).$$

The gradient of  $\varphi_x$  can be expressed in terms of the original data of the problem

$$(1.6) \quad g^k = \nabla_J f(x^k) - D^T (B^{-1})^T \nabla_I f(x^k)$$

and is called the reduced gradient for the partition  $N = I \oplus J$ ; the stepsize  $t_k$  must be selected by a search along the arc of curve starting from  $x^k$  given by

$$(1.7) \quad x(t) = (\psi(x_J^k + t p^k), x_J^k + t p^k)$$

to produce a sufficient decrease of the objective function. With this framework, Gabay and Luenberger proposed idealized methods based on the reduced gradient extending the steepest descent method, Newton's method and a Quasi-Newton method and analyzed their convergence properties. Notice that if  $x^* \notin U_x^k \cap C$  it is still possible to find a new approximation  $x^{k+1}$  according to (1.5) by restricting the search along the curve (1.7) to an interval  $[0, \tilde{t}]$  such that  $x^{k+1} \in U_x^k \cap C$ ; at  $x^{k+1}$  a new partition  $N = I \oplus J$  can be constructed and an iterate  $x^k$  is eventually reached such that  $x^* \in U_x^k \cap C$ .

This framework exploits the fact that the constrained set  $C$  is locally diffeomorphic to an open set  $V$  of the Euclidean space  $\mathbb{R}^{n-m}$  and that the corresponding pieces of  $\mathbb{R}^{n-m}$  can be "glued together" by diffeomorphisms. This informal description characterizes  $C$  as a differential submanifold of  $\mathbb{R}^n$ . The  $n - m$  coordinates  $z$  are said to form a local coordinate system of the manifold  $C$  around  $x$ .

In this paper we generalize the approach of (Ref. 1) by defining directly descent methods for the minimization of the real-valued function  $f$  over the differential manifold  $C$ , independently of the choice of coordinate systems. In order to define the gradient (field) of  $f$  on  $C$  we must endow  $C$  with a Riemannian structure (we similarly use implicitly the Euclidean structure of  $\mathbb{R}^n$  to define gradient-related methods for unconstrained minimization). Such methods consist in generating from an approximation  $x^k$  a new iterate  $x^{k+1}$  on the geodesic

curve starting from  $x^k$  and tangent to a direction defined by the gradient of  $f$  on  $C$  at  $x^k$ . This approach was inaugurated by Luenberger (Ref. 2) to study the convergence of Rosen's Gradient Projection method (Ref. 3); in (Ref. 2) the Riemannian structure on  $C$  is the one induced by the Euclidean structure of  $\mathbb{R}^n$ . In order to define a family of methods which includes both the Gradient Projection and the Reduced Gradient method we define in this paper descent methods for a general Riemannian structure on  $C$ ; our approach is inspired by the recent work of Lichnerowski (Ref. 4) where a gradient and a conjugate-gradient methods along geodesics are analyzed.

In Section 2 we precisely show that, under the regularity assumption,  $C = c^{-1}(0)$  is a differential manifold and introduce a nonlinear change of coordinates in  $\mathbb{R}^n$  which is useful for the practical implementation of our methods. We also give an estimate of the diameter of the neighborhood  $U_x \cap C$  on which a coordinate system around  $x$  can be defined. Section 3 is of a tutorial nature and presents important results on the Geometry of a Riemannian manifold which are used in our analysis. Section 4 constitutes the core of the paper. We first establish optimality conditions in term of the restriction of  $f$  to the manifold  $C$  and show their relations with the traditional Lagrange multipliers rules. We then define the gradient field of  $f$  on  $C$  and show how it depends upon the Riemannian metric; we thus defined the reduced gradient in a local coordinate system which yields Abadie's reduced gradient used in (Ref. 1) and Rosen's projected gradient as special cases. We then define the steepest descent method, the Newton's method and Quasi-Newton methods along geodesics and analyse their convergence properties. In section 5 these methods are specified in the framework of a coordinate system. The analysis of section 2 then provides a practical scheme for the implementation of the algorithms which can be interpreted as the sequential tangent-restoration approach introduced by Rosen (Ref. 3) and used by Abadie, Guigou (Ref. 5), Miele and al. (Ref. 6). We give in particular two efficient versions of a Quasi-Newton method for constrained minimization based on a Generalized Broyden-Fletcher-Goldfarb-Shanno update formula. Finally in Section 6 we indicate how our approach can be used to handle problems with inequality constraints.



## 2 - THE CONSTRAINED SET AS A MANIFOLD

In this paragraph we study the properties of the constrained set  $C = c^{-1}(0)$ , defined by the map  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  ( $m \leq n$ ) assumed to be  $\mathcal{C}^\sigma$  differentiable ( $\sigma \geq 2$ ). Recall that 0 is a regular value of  $c$  iff  $c'(x) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$  is of full rank  $m$  for all  $x \in C$ .

**THEOREM 2.1** Let  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  ( $m \leq n$ ) be a  $\mathcal{C}^\sigma$  map such that 0 is a regular value of  $c$ . Then the set  $C = c^{-1}(0)$  is an  $(n - m)$  dimensional submanifold of  $\mathbb{R}^n$  of class  $\mathcal{C}^{\sigma-1}$ . ■

Proof : To prove this result, often referred to as the regular value theorem, we follow Milnor (Ref. 7) and introduce a nonlinear change of coordinates in  $\mathbb{R}^n$  which will be useful for our analysis.

Fix a point  $x \in C$ ; thus  $c(x) = 0$ . Since 0 is a regular value of  $c$ , the derivative  $c'(x) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$  must map  $\mathbb{R}^n$  onto  $\mathbb{R}^m$ . The nullspace  $\mathcal{N}[c'(x)]$  is therefore an  $(n - m)$  dimensional subspace of  $\mathbb{R}^n$ . Now choose a linear map  $Z_x \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^{n-m})$  that is nonsingular on this subspace, i.e. such that.

$$(2.1) \quad \mathcal{N}[Z_x] \cap \mathcal{N}[c'(x)] = \{0\}.$$

Define the mapping  $s_x : \mathbb{R}^n \rightarrow \mathbb{R}^m \times \mathbb{R}^{n-m}$  by

$$(2.2) \quad s_x(y) = \begin{pmatrix} c(y) \\ Z_x(y-x) \end{pmatrix} \quad \text{for all } y \in \mathbb{R}^n.$$

Notice that  $s_x(x) = 0$  and that the derivative  $s'_x(x) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^n)$ , given by the formula

$$(2.3) \quad s'_x(x)(y) = \begin{pmatrix} c'(x)(y) \\ Z_x(y) \end{pmatrix} \quad \text{for all } y \in \mathbb{R}^n,$$

is non singular by (2.1).

The inverse function theorem implies that  $s_x$  maps some neighborhood  $U_x$  of  $x$  in  $\mathbb{R}^n$  diffeomorphically onto a neighborhood  $W \times V$  of 0 in  $\mathbb{R}^m \times \mathbb{R}^{n-m}$ . Under the change of coordinates  $s_x$ , the set  $C$  is transformed locally in the  $(n-m)$  dimensional subspace  $\{0\} \times \mathbb{R}^{n-m}$  of  $\mathbb{R}^n$  since  $s_x$  maps  $U_x \cap C$  diffeomorphically onto  $\{0\} \times V$ . The map  $Z_x$  defines by restriction a coordinate system of  $C$  around  $x$  on the coordinate domain  $U_x \cap C$  and the components  $(z_1, \dots, z_{n-m})$  of the image  $z = Z_x(y-x)$  of  $y \in U_x \cap C$  are called the local coordinates of  $y$ . The inverse diffeomorphism  $s_x^{-1}$  allows to define a

local (non linear) parametrization  $\theta_x : V \rightarrow U_x \cap C$  of  $C$  around  $x$  by

$$(2.4) \quad y = \theta_x(z) = s_x^{-1}(0, z).$$

Thus  $C$  is an  $(n-m)$  dimensional manifold. Provided the map  $Z_x$  is chosen as a  $\mathcal{C}^{\sigma-1}$  differentiable function of  $x$ , the change from local coordinates around  $x$  to local coordinates around  $x'$  (such that  $U_x \cap U_{x'} \cap C \neq \emptyset$ ), defined by

$$Z_{x'} \circ \theta_x : Z_x(U_x \cap U_{x'} \cap C) \rightarrow Z_{x'}(U_x \cap U_{x'} \cap C),$$

is  $\mathcal{C}^{\sigma-1}$  differentiable and the manifold  $C$  is of class  $\mathcal{C}^{\sigma-1}$ . ■

Since we are mainly interested by computational procedures we may wish to specify further the choice of the map  $Z_x$  defining the coordinate system of  $C$ . Let us denote by  $A_x$  the  $m \times n$  Jacobian matrix of  $c$  at  $x$ , while we keep  $Z_x$  to denote the  $(n-m) \times n$  matrix representing the linear map  $Z_x$ . The Jacobian matrix of the map  $s_x$  is the  $n \times n$  matrix  $S_x$  given in partitioned form by

$$(2.5) \quad S_x = \begin{bmatrix} A_x \\ Z_x \end{bmatrix}.$$

Notice that the regularity assumption implies that  $\text{rank}(A_x) = m$  for all  $x \in C$ . Recall that a right inverse  $M^-$  for a  $l \times n$  matrix  $M$  of full rank  $l \leq n$  is a  $n \times l$  matrix of full rank  $l$  such that

$$M M^- = I.$$

Such an inverse exists but is not unique. The following result relates the choice of  $Z_x$  (and  $Z_x^-$ ) to a particular choice of a right inverse  $A_x^-$  of  $A_x$  in order to express the inverse  $S_x^{-1}$  in a simple form.

**PROPOSITION 2.1** Let  $A_x^-$  be a right inverse for  $A_x$ . Then there exists an  $(n-m) \times n$  matrix  $Z_x$  of full rank  $(n-m)$  and a right inverse  $Z_x^-$  satisfying

$$(2.6) \quad Z_x \cdot A_x^- = 0, \quad A_x \cdot Z_x^- = 0$$

such that the matrix  $S_x$  given by (2.5) is non-singular and its inverse is given by

$$(2.7) \quad S_x^{-1} = [A_x^-, Z_x^-] \quad \blacksquare$$

Proof : A necessary and sufficient condition for the non-singularity of  $S_x$  is

$$(2.8) \quad \mathcal{N}(Z_x) \cap \mathcal{N}(A_x) = \{0\}.$$

Since  $A_x$  is of full rank  $m$ ,  $\mathcal{N}(A_x)$  is a subspace of  $\mathbb{R}^n$  of dimension  $n - m$ . It follows from (2.8) that  $\mathcal{N}(Z_x)$  must be a complement of  $\mathcal{N}(A_x)$  in  $\mathbb{R}^n$ , i.e. an  $m$ -dimensional subspace; hence  $Z_x$  must be of full rank  $(n-m)$ . The right inverse  $A_x^-$  being of rank  $m$ , its columns span a subspace  $\mathcal{R}(A_x^-)$  of  $\mathbb{R}^n$  of dimension  $m$  such that

$$\mathcal{N}(A_x) \cap \mathcal{R}(A_x^-) = \{0\}$$

by definition of a right inverse. Hence  $\mathcal{N}(Z_x) = \mathcal{R}(A_x^-)$ , yielding

$$Z_x \cdot A_x^- = 0.$$

A similar argument shows that  $\mathcal{R}(Z_x^-) = \mathcal{N}(A_x)$ , i.e.  $A_x \cdot Z_x^- = 0$ .

Formula (2.7) is established by observing that

$$y = A_x^- a + Z_x^- z$$

satisfies the equation

$$S_x y = \begin{pmatrix} a \\ z \end{pmatrix}$$

for all  $a \in \mathbb{R}^m$ ,  $z \in \mathbb{R}^{n-m}$  and is the unique solution since  $S_x$  is non singular. ■

Define the  $n \times n$  matrix  $P_x$  by

$$(2.9) \quad P_x = I - A_x^- A_x.$$

It is the matrix of a projection onto  $\mathcal{N}(A_x)$  since  $A_x P_x = 0$  and the matrices  $Z_x$  and  $Z_x^-$  defined by Proposition 2.1 satisfy

$$(2.10) \quad Z_x^- \cdot Z_x = P_x.$$

Example 2.1 : (Partitioned Right Inverse). After possibly a permutation of columns, the matrix  $A_x$  can be partitioned into

$$(2.11) \quad A_x = [B, D],$$

such that B is an  $m \times m$  non-singular matrix. It is easy to verify that

$$(2.12) \quad A_x^- = \begin{bmatrix} B^{-1} \\ 0 \end{bmatrix}$$

is a right inverse for  $A_x$  and that

$$(2.13) \quad Z_x = [0, I_{n-m}] \quad , \quad Z_x^- = \begin{bmatrix} -B^{-1}D \\ I_{n-m} \end{bmatrix}$$

satisfy proposition 2.1. In this case the local coordinate system around  $x$  is formed by some  $(n-m)$  coordinates in  $\mathbb{R}^{n-m}$ ; it is the system used by Gabay-Luenberger (Ref. 1).

From a computational view point this choice requires the inversion of the  $m \times m$  matrix B which can be performed by Gaussian elimination with partial pivoting in the order of  $m^3/3$  multiplications; but the partition selected may produce a ill-conditioned matrix B even though  $A_x$  is not ill-conditioned. ■

Example 2.2 : (QR pseudo-inverse). A classical choice for a right inverse of  $A_x$  is the Penrose pseudo-inverse  $A_x^+$ , which in the case of a full rank matrix  $A_x$  can be expressed as

$$(2.14) \quad A_x^+ = A_x^T (A_x A_x^T)^{-1}$$

Notice that the matrix  $P_x$  defined by (2.9) is now

$$P_x = I - A_x^T (A_x A_x^T)^{-1} A_x,$$

the orthogonal projector onto  $\mathcal{N}(A_x)$  in the Euclidian space  $\mathbb{R}^n$ . Relations (2.6) are satisfied by choosing  $Z_x$  such that  $\mathcal{N}(Z_x)$  is the orthogonal complement of  $\mathcal{N}(A_x)$  in  $\mathbb{R}^n$  and taking for  $Z_x^-$  its pseudo-inverse. An equivalent characterization of  $Z_x$  is

$$(2.15) \quad \mathcal{R}(Z_x^T) = \mathcal{N}(A_x)$$

and a particularly convenient choice satisfying (2.15) consists in taking for columns of  $Z_x^T$  an orthonormal-basis of  $\mathcal{N}(A_x)$ ; in this case notice that  $Z_x^- = Z_x^T$  and hence

$$(2.16) \quad Z_x^{-T} \cdot Z_x^- = I_{n-m}.$$

From a computational viewpoint, notice that the formula (2.14) involves the inverse of a matrix of condition number  $[K(A)]^2$  and should be avoided in practice. Besides, while it is theoretically possible to generate the columns of  $Z_x^T$  by Gram-Schmidt orthogonalization method, such a procedure is numerically unstable. It is however possible to obtain  $A_x^+$  and  $Z_x$  in a numerically efficient way once  $A_x$  is factorized as

$$(2.17) \quad A_x = \Lambda Q,$$

where  $Q$  is an  $(n \times n)$  orthogonal matrix ( $Q Q^T = Q^T Q = I$ ) and  $\Lambda$  an  $(m \times n)$  matrix of the form

$$\Lambda = [L, 0]$$

with  $L$   $(m \times m)$  lower triangular. This "Q-R decomposition" can be carried out in a numerically stable way using Householder's transformations with approximately  $(n - m/3) m^2$  multiplications (see e.g. Stewart (Ref. 8)).

Partition  $Q$  into

$$(2.18) \quad Q = \begin{bmatrix} Q_1 \\ Q_2 \end{bmatrix}$$

where  $Q_1$  and  $Q_2$  are respectively  $m \times n$  and  $(n - m) \times n$  submatrices of  $Q$ . The pseudo-inverse (2.14) can now be efficiently computed according to

$$(2.19) \quad A_x^+ = Q_1^T L^{-1},$$

while the columns of  $Q_2^T$  form an orthonormal basis of  $\mathcal{N}(A_x)$ ; hence

$$(2.20) \quad Z_x = Q_2, \quad Z_x^- = Q_2^T$$

satisfy the conditions of Proposition 2.1. Notice again that  $Z_x^{-T} Z_x^- = I$ . ■

Remark 2.1. The singular value decomposition is commonly used to compute the pseudo-inverse of a matrix but requires more computational effort than the Q-R decomposition (see Steward (Ref. 8)). It is useful however when the matrix is not of full rank a situation that we have excluded by hypothesis. ■

Remark 2.2. We must observe that if the map  $c$  is a  $\mathcal{C}^\sigma$  differentiable, the map  $Z_x$  satisfying Proposition 2.1 can only be chosen  $\mathcal{C}^{\sigma-1}$  differentiable with respect to  $x$ . The resulting coordinate systems on  $C$  give only a  $\mathcal{C}^{\sigma-1}$  differential structure to the submanifold. Since  $\sigma \geq 2$ ,  $C$  will still be a differential manifold. ■

We conclude this subsection by giving an estimate of the coordinate domain  $U_x \cap C$  in terms of the map  $c$ . Take  $U_x \subset B(x, \alpha)$ , the open ball in  $\mathbb{R}^n$  of centre  $x$  and radius  $\alpha > 0$ ; let  $\beta = \|A_x^-\|$ ,  $\xi = \|Z_x^-\|$ ,

$$\gamma_i = \text{Max}_{\tilde{x} \in B(x, \alpha)} \text{Max}_{\|y\|=1} |c''_i(\tilde{x})(y)(y)|,$$

where  $\|\cdot\|$  stands for the  $\ell_2$  norm; let

$$\gamma = \left( \sum_{i=1}^n \gamma_i \right)^{1/2}.$$

THEOREM 2.2. The local parametrization  $\theta_x$  defined by (2.4) around  $x$  maps diffeomorphically the neighborhood  $B(0, 1/(2\gamma\xi))$  of the origin in  $\mathbb{R}^{n-m}$  onto a neighborhood  $U_x \cap C$  of  $x$  in  $C$  where  $U_x \subset B(x, 1/(\beta\gamma))$ . ■

Proof : The map  $s_x$  given by (2.2) defines, by restriction, a diffeomorphism of  $U_x \cap C$  onto a neighborhood  $V$  of the origin in  $\mathbb{R}^{n-m}$ . To get an estimate of  $U_x$ , we look for an estimate of  $V$  such that given  $z \in V$ , there exists a unique  $y \in C$  such that  $s_x(y) = (0, z)$ .

Let  $\|z\| = \rho$ . Proposition 2.1 shows that we can look for

$$(2.21) \quad y = x + Z_x^- z + A_x^- w,$$

with  $w \in \mathbb{R}^m$  such that

$$(2.22) \quad h(w) = c(x + Z_x^- z + A_x^- w) = 0.$$

We shall show that if  $\rho \leq 1/(2 \beta \gamma \xi)$  there exists a unique solution of (2.22) in  $B(0, 1/(2 \beta^2 \gamma))$  by a constructive proof.

Consider the iteration

$$(2.23) \quad w^{k+1} = w^k - h(w^k) \quad k = 0, 1, \dots$$

starting from  $w^0 = 0 \in \mathbb{R}^m$ . We construct a majorizing sequence  $\{t_k\}$  of non negative real numbers, i.e. such that

$$(2.24) \quad \|w^k\| \leq t_k \quad \text{for all } k.$$

Applying Taylor formula to the map  $c$  around  $x$  we obtain

$$(2.25) \quad w^{k+1} = - \int_0^1 (1-t) c''(x + t Z_x^- z + t A_x^- w^k) (Z_x^- z + A_x^- w^k) (Z_x^- z + A_x^- w^k) dt;$$

define the sequence  $\{t_k\}$  by the iteration

$$(2.26) \quad t_{k+1} = \frac{\gamma}{2} (\beta t_k + \xi \rho)^2 \quad k = 0, 1, \dots,$$

starting from  $t_0 = 0$ . If  $\|w^k\| \leq t^k$ , (2.25) shows that  $\|w^{k+1}\| \leq t^{k+1}$ , and by induction  $\{t_k\}$  is actually a majorizing sequence. It is easy to show that if

$$(2.27) \quad \rho \leq \frac{1}{2 \beta \gamma \xi},$$

the sequence  $\{t_k\}$  is monotonically increasing and converges to

$$(2.28) \quad t_\rho^* = \frac{1 - \beta \gamma \xi \rho - (1 - 2 \beta \gamma \xi \rho)^{1/2}}{\beta^2 \gamma} > 0;$$

hence,  $w^k$  remains in  $B(0, t_\rho^*)$  for all  $k$ .

Consider now  $w^{k+1} - w^k$ , which can be written

$$(2.29) \quad \begin{aligned} w^{k+1} - w^k = & - [h(w^k) - h(w^{k-1}) - h'(w^{k-1})(w^k - w^{k-1})] \\ & - [h'(w^{k-1}) - h'(0)](w^k - w^{k-1}) + [I - h'(0)](w^k - w^{k-1}). \end{aligned}$$

Assuming the induction hypothesis (trivially satisfied for  $k = 1$ )

$$(2.30) \quad \|w^k - w^{k-1}\| \leq t_k - t_{k-1},$$

we obtain, with the use of Taylor formula, the estimate

$$(2.31) \quad \begin{aligned} \|w^{k+1} - w^k\| &\leq \frac{\gamma}{2} [\beta^2 (t_k - t_{k-1})^2 + 2\beta^2 t_{k-1} (t_k - t_{k-1}) \\ &\quad + 2\beta\gamma\rho (t_k - t_{k-1})] \\ &= \frac{\gamma}{2} [(\beta t_k + \xi\rho)^2 - (\beta t_{k-1} + \xi\rho)^2] = t_{k+1} - t_k, \end{aligned}$$

which proves that (2.30) holds for all  $k$  and that  $\|w^{k+1} - w^k\| \rightarrow 0$  since  $t_k \nearrow t_\rho^*$ ; hence the convergence of the sequence  $\{w^k\}$  defined by (2.23) to  $w$  such that  $\|w\| \leq t_\rho^*$ , unique solution of the equation (2.22) in  $B(0, t_\rho^{**})$  where  $t_\rho^{**}$  is the other fixed point of (2.26) given by

$$(2.32) \quad t_\rho^{**} = \frac{1 - \beta\gamma\xi\rho + (1 - 2\beta\gamma\xi\rho)^{1/2}}{\beta^2\gamma} \geq t_\rho^*.$$

Notice that for the boundary value  $\rho = 1/(2\beta\gamma\xi)$ ,  $t_\rho^{**} = t_\rho^* = 1/(2\beta^2\gamma)$ .

Then for any  $z \in B(0, 1/(2\beta\gamma\xi))$  there exists a unique solution  $w$  of (2.22) in  $B(0, 1/(2\beta^2\gamma))$  and  $\|y - x\| = \|Z_x^- z + A_x^- w\| \leq 1/\beta\gamma$ . ■

Remark 2.3. The iteration (2.23) can be interpreted as a secant(or modified Newton) method for solving the equation (2.21) where the Jacobian  $h'(w^k)$  is approximated by the identity matrix. ■



### 3 - THE GEOMETRY OF THE CONSTRAINED MANIFOLD

#### 3.1. - A Riemannian structure on C

Given a point  $x \in C$  we define the tangent space  $T_x$  to the manifold  $C$  at  $x$  as the  $(n-m)$  dimensional subspace of  $\mathbb{R}^n$

$$(3.1) \quad T_x = \mathcal{N}(A_x) = \{v \in \mathbb{R}^n \mid A_x v = 0\}.$$

If we choose the local coordinate system defined in the neighborhood  $U_x \cap C$  of  $x$  in  $C$  by

$$(3.2) \quad z(y) = Z_x(y - x) \quad \text{for } y \in U_x \cap C,$$

we have, as noticed in the proof of Proposition 2.1.,  $\mathcal{N}(A_x) = \mathcal{R}(Z_x^-)$  and

$$(3.3) \quad T_x = \{Z_x^- p \mid p \in \mathbb{R}^{n-m}\} \subset \mathbb{R}^n.$$

The equivalence of (3.3) with (3.1) shows that this last definition is independent of the coordinate system (see e.g. Guillemin, Pollack (Ref. 9))

It is convenient to endow  $T_x$  with a positive bilinear form  $\gamma_x : T_x \times T_x \rightarrow \mathbb{R}$ , called the Riemannian metric. The form  $\gamma_x$  is a smooth function of  $x$  and defines a Riemannian structure  $\gamma$  on  $C$ . A natural choice consists in taking

$$(3.4) \quad \gamma_x^E(v, w) = \langle v, w \rangle_n \quad \forall v, w \in T_x \subset \mathbb{R}^n$$

where  $\langle \cdot, \cdot \rangle_n$  denotes the ordinary scalar product on  $\mathbb{R}^n$ ;  $\gamma_x^E$  is called the Riemannian metric on  $C$  induced by the Euclidean structure of  $\mathbb{R}^n$ . In the local coordinate system (3.2), the induced metric can be expressed as

$$\gamma_x^E(Z_x^- p, Z_x^- q) = \langle p, Z_x^{-T} \cdot Z_x^- q \rangle_{n-m} \quad \forall p, q \in \mathbb{R}^{n-m},$$

which coincides with the ordinary scalar product on  $\mathbb{R}^{n-m}$  iff  $Z_x^{-T} \cdot Z_x^- = I$  (see example 2.2. for such a case). It may be convenient to directly define the form in the local coordinate system by

$$(3.5) \quad \gamma_x^Z(Z_x^- p, Z_x^- q) = \langle p, q \rangle_{n-m} \quad \forall p, q \in \mathbb{R}^{n-m}.$$

Notice that  $\gamma_x^Z$  can be viewed as the Riemannian metric on  $C$  induced by a scalar product in  $\mathbb{R}^n$  defined locally around  $x$  by

$$\langle v, w \rangle_{S_x} = \langle S_x v, S_x w \rangle_n \quad \forall v, w \in \mathbb{R}^n,$$

where  $S_x$  is the matrix defined by (2.5). A general Riemannian structure on  $C$  can be defined in a local coordinate system around  $x$  by

$$(3.6) \quad \gamma_x^G(Z_x^{-1} p, Z_x^{-1} q) = \langle p, G_x q \rangle_{n-m},$$

where  $G_x = (g_{ij})$  is a positive definite symmetric matrix of dimension  $n - m$ , smooth function of  $x$ . (The terminology smooth is used here and in the following for continuously differentiable to a sufficiently large order for all formulae to have sense).

### 3.2. - Covariant Differentiation and Geodesics

We define a vector field  $V$  on the submanifold  $C$  of  $\mathbb{R}^n$  as a smooth map  $V : C \rightarrow \mathbb{R}^n$  such that  $V(x) \in T_x$  for all  $x$ .

Let  $x \in C$ . Given a vector  $v \in T_x$  and a vector field  $W$  on  $C$ , we define a new vector  $D_v W \in T_x$ , called the covariant derivative of  $W$  along  $v$ ; the application  $\tau_x(W) : T_x \rightarrow T_x$  defined by

$$(3.7) \quad \tau_x(W)(v) = D_v W$$

must satisfy

$$(3.8a) \quad \tau_x(W)(\alpha_1 v_1 + \alpha_2 v_2) = \alpha_1 \tau_x(W)(v_1) + \alpha_2 \tau_x(W)(v_2)$$

$$(3.8b) \quad \tau_x(fW)(v) = f(x) \tau_x(W)(v) + f'(x)(v) W(x)$$

where  $f$  is any real-valued smooth function on  $C$ , and specifies an affine connexion on  $C$  at  $x$ . Let now  $V$  and  $W$  be vector fields on  $C$ ; we define the field  $D_V W$ , the covariant derivative of  $W$  with respect to  $V$  on  $C$ , by its values

$$(3.9) \quad D_V W(x) = D_v(W) \quad \text{where } v = V(x) \in T_x$$

The affine connexion is thus specified globally on  $C$  (see Milnor (Ref. 10)).

Given the local coordinate system (3.2) around  $x$ , the column vectors  $e_i$  ( $i = 1, \dots, n-m$ ) of the matrix  $Z_x^{-1}$  form a basis of  $T_x$  in  $\mathbb{R}^n$ . The  $\mathcal{C}^{\sigma-1}$  maps  $E_i : C \rightarrow \mathbb{R}^n$  such that  $E_i(x) = e_i$  are vector fields and are said to form the associated base fields around  $x$ . Properties (3.8) show that the affine connexion on  $C$  is determined by the fields  $D_{E_i} E_j$ . It is customary to express these vector fields in the base fields as

$$(3.10) \quad D_{E_i} E_j = \sum_{k=1}^{n-m} \Gamma_{ij}^k E_k.$$

The  $(n-m)^3$  real-valued smooth functions  $\Gamma_{ij}^k$  determine the connexion around  $x$  and are called the coefficients of the connexion.

A parametrized curve in  $C$  is a smooth map  $x(\cdot)$  from the real numbers into  $C \subset \mathbb{R}^n$ . The velocity vector field  $\dot{x}$  is defined by

$$(3.11) \quad \dot{x}(t) = \frac{dx}{dt}(t) \in T_{x(t)} \quad \text{for all } t \in \mathbb{R}$$

A vector field  $V$  on  $C$  defines, by restriction, a vector field  $v(\cdot)$  along the curve  $x(\cdot)$  which assigns to each  $t \in \mathbb{R}$  a tangent vector

$$(3.12) \quad v(t) = V(x(t)) \in T_{x(t)}.$$

The vector field  $v(\cdot)$  is said to be a parallel vector field along the curve  $x(\cdot)$  if its covariant derivative, denoted  $\frac{Dv}{dt}$ , is identically zero, i.e.

$$(3.13) \quad \frac{Dv}{dt} = D_{\dot{x}} V = 0.$$

Using the local coordinate system (3.2) around a point  $x$  of the curve and the associated base fields  $E_i$ , the vector field  $v(\cdot)$  can be uniquely expressed as

$v = \sum_{i=1}^{n-m} v_i E_i$  where  $v_i$  are smooth real-valued functions on  $\mathbb{R}$ , while the velocity

field  $\dot{x} = \sum_{i=1}^{n-m} \frac{dz_i}{dt} E_i$  where  $z_i(t) = z[x(t)]$ . It follows from (3.7) (3.8) (3.9)

(3.10) (3.13) that the functions  $v_i$  must satisfy the system of linear differential equations

$$(3.14) \quad \frac{dv_k}{dt} + \sum_{i,j=1}^{n-m} \frac{dz_i}{dt} \Gamma_{ij}^k v_j = 0 \quad k = 1, 2, \dots, n-m$$

and we have the following existence and uniqueness result (see also Hicks (Ref. 11)).

PROPOSITION 3.1. Let  $x(\cdot)$  be a parametrized curve in  $C$  defined on  $[0, T]$ . For each vector  $v$  in  $T_{x(o)}$  there is a unique parallel vector field  $v(\cdot)$  along  $x(\cdot)$  such that  $v(o) = v$ . The map  $\pi_o^t : T_{x(o)} \rightarrow T_{x(t)}$  defined by  $\pi_o^t(v(o)) = v(t)$  is a linear isomorphism called parallel translation along  $x(\cdot)$  from  $x(o)$  to  $x(t)$ . ■

The following result is a fundamental tool in Riemannian geometry (see Milnor (Ref. 10)).

PROPOSITION 3.2. There exists a unique (symmetric) connexion on a Riemannian manifold such that parallel translation preserves the Riemannian metric. ■

In a system of local coordinates  $\{z_i\}$  the coefficient functions  $\Gamma_{ij}^k$  defining the connexion are uniquely determined by the symmetric matrix function  $G$  defining the Riemannian metric in (3.6) and satisfy for  $i, j, k = 1, 2, \dots, n-m$

$$(3.15) \quad \Gamma_{ij}^k = \sum_{\ell} \frac{1}{2} \left( \frac{\partial g_{jk}}{\partial z_i} + \frac{\partial g_{ik}}{\partial z_j} - \frac{\partial g_{ij}}{\partial z_k} \right) (G^{-1})_{\ell k},$$

called the second Christoffel identity. Notice that  $\Gamma_{ij}^k = \Gamma_{ji}^k$  for all  $i, j$ , which indicates that the connexion is symmetric. Notice also that if  $G(x)$  remains constant (for instance if  $G(x) = I$ ),  $\Gamma_{ij}^k \equiv 0$ .

A parametrized curve  $x$  of  $C$  is called a geodesic if its velocity field  $\dot{x}$  is parallel along  $x$ , i.e.  $\frac{D\dot{x}}{dt} = 0$ . In terms of the local coordinate system (3.2) the local coordinate functions  $z = (z_k)$   $k = 1, \dots, n-m$ , must satisfy the system of second order differential equations

$$(3.16) \quad \frac{d^2 z_k}{dt^2} + \sum_{i,j} \Gamma_{ij}^k(z) \frac{dz_i}{dt} \frac{dz_j}{dt} = 0,$$

derived from (3.14) with  $v_k = \frac{dz_k}{dt}$ . From now on we consider  $C$  endowed with a Riemannian metric and work with the unique connexion compatible with it, i.e. the coefficients of which satisfy the second Christoffel identity. The following result about the solutions of (3.16) establishes the equivalence of our definition of geodesics with the elementary one (curves of minimal lengths).

PROPOSITION 3.3. Let  $W$  be a connected compact subset of  $C$  in  $\mathbb{R}^n$ . Given  $x \in W$  and  $p \in T_x$ , there exists a unique geodesic curve  $x(\cdot)$  such that  $x(o) = x$ ,  $\dot{x}(o) = p$ ; the map  $x$  is either defined for all  $t \in \mathbb{R}$  and takes its value in  $W$ , or defined on the interval  $[-T, T]$  and  $x(T)$  (or  $x(-T)$ ) belongs to the boundary of  $W$ . Any two points of  $W$  can be joined by a unique geodesic which minimizes the arc length between the points. ■

See Milnor (Ref. 10) or Bishop, Crittenden (Ref. 11) for a proof. The geodesic curve will be denoted

$$(3.17) \quad x(t) = \exp_x(tp)$$

and the local coordinate system  $\exp_x^{-1} : U_x \cap C \rightarrow V \subset T_x = \mathbb{R}^{n-m}$  called normal coordinates around  $x$ . Observe that in such coordinates the geodesic is locally the parametrized line interval  $\{tp \mid t \in (-\epsilon, +\epsilon)\}$ . Finally notice that if in a local coordinate system the Riemannian metric is defined by a constant metric  $G$ , the local coordinate functions  $z = (z_k) \ k = 1, \dots, n-m$  defining the geodesic around  $x$  satisfy the second order differential equation

$$\frac{d^2 z_k}{dt^2} = 0 \quad i = 1, 2, \dots, n-m$$

since  $\Gamma_{ij}^k = 0$ ; hence  $z(t) = tp$  and the local coordinate system is normal.

#### 4 - DESCENT METHODS ALONG GEODESICS

##### 4.1. - Optimality conditions

We turn now to the solution of the nonlinear programming problem (1.1) which we now write

$$(4.1) \quad \text{Min } \{f(x) \mid x \in C\},$$

i.e. the problem of minimizing the differentiable real-valued function  $f$  on the differential manifold  $C$  endowed with a smooth Riemannian structure  $\gamma$ .

Given  $x \in C$  we define the derivative on  $C$  of  $f$  at  $x$  as the linear form on  $T_x$ ,  $D_C f(x) : T_x \subset \mathbb{R}^n \rightarrow \mathbb{R}$ , such that

$$(4.2) \quad D_C f(x)(v) = f'(x)(v) \quad \text{for all } v \in T_x,$$

where  $f'(x)$  is the ordinary derivative at  $x$  of  $f$  considered as a function on  $\mathbb{R}^n$ .

We can now state a first-order necessary optimality condition for problem (4.1).

**THEOREM 4.1** Let  $x^* \in C$  be a local minimum of  $f$  on the Riemannian manifold  $C$ . Then

$$(4.3) \quad D_C f(x^*) = 0. \quad \blacksquare$$

Proof : Since  $x^*$  is a local minimum of  $f$  on  $C$

$$(4.4) \quad f(x) \geq f(x^*) \quad \forall x \in U_{x^*} \cap C,$$

where  $U_{x^*}$  is a neighborhood of  $x^*$  in  $\mathbb{R}^n$ . According to Proposition 3.3 (4.4) is equivalent to

$$(4.5) \quad f(x(t)) \geq f(x^*)$$

for any geodesic curve  $x(\cdot)$  starting from  $x^*$  and all  $t \in (-\epsilon, +\epsilon)$ . Applying the mean value theorem to the function  $f \circ x$  between 0 and  $t$ , (4.5) yields

$$(4.6) \quad f'(x^*) \dot{x}(0) = f'(x^*) \cdot v \geq 0 \quad \forall v \in T_{x^*};$$

hence

$$D_C f(x^*)(v) = 0 \quad \forall v \in T_{x^*}$$

by applying (4.6) to both  $v$  and  $-v$ .  $\blacksquare$

Conversely, every point  $x^* \in C$  such that  $D_C f(x^*) = 0$  is called a critical point of  $f$ .

Given  $v \in T_{x^*}$  let  $V$  be a smooth vector field on  $C$  such that  $V(x) = v$ . Define the  $\mathcal{C}^{\sigma-1}$  differentiable real-valued function  $Vf : C \rightarrow \mathbb{R}$  by

$$(4.7) \quad Vf(x) = D_C f(x)(v).$$

Since  $\sigma \geq 2$  this function has a derivative on  $C$  at  $x$ ,  $D_C Vf(x) \in L(T_x, \mathbb{R})$

If  $x^*$  is a critical point, we define the Hessian of  $f$  on  $C$ , as the bilinear form  $Hf(x^*) : T_{x^*} \times T_{x^*} \rightarrow \mathbb{R}$  given by

$$(4.8) \quad Hf(x^*)(v, w) = D_C Vf(x^*)(w) \quad \forall v, w \in T_{x^*}$$

This form is well-defined (i.e. independent of the choice of the vector field  $V$ ) and symmetric (See Milnor (Ref. 10)). This definition holds only at a critical point.

A critical point  $x^*$  is said to be non-degenerated iff  $Hf(x^*)$  is non-degenerated, i.e.  $Hf(x^*)(v,v) = 0$  implies  $v = 0$ . In the following we assume that all the critical points of  $f$  on  $C$  are non-degenerated;  $f$  is called a Morse function on  $C$  (see Milnor (Ref. 10)). It follows that the critical points are isolated. We can then give a second-order optimality condition in stronger terms than usually expressed (see e.g. Luenberger (Ref. 13)).

**THEOREM 4.2** A point  $x^* \in C$  is a strict local minimum of a Morse function  $f$  on the Riemannian manifold  $C$  if and only if  $x^*$  is a critical point of  $f$  and the Hessian form  $Hf(x^*)$  is positive definite. ■

Proof : Let  $x(\cdot)$  be the geodesic curve starting from  $x^*$  with tangent  $v$  and  $V$  the unique parallel vector field along  $x(\cdot)$  such that  $V(x^*) = v$ . Applying Taylor formula to the  $\mathcal{C}^\sigma$  differentiable function  $\text{fo}x : \mathbb{R} \rightarrow \mathbb{R}$  yields

$$(4.9) \quad f[x(t)] - f(x^*) = D_C f(x^*)(v)t + \frac{1}{2} D_C V f(x(\theta t)) (V(x(t)))t^2$$

with  $\theta \in (0,1)$ . If  $x^*$  is a local minimum, the first member of the right hand side vanishes by Theorem 4.1 while the second member must be non negative; by continuity  $Hf(x^*)(v,v) \geq 0 \forall v \in T_{x^*}$  and since  $f$  is a Morse function strict inequality holds for all  $v \neq 0$  and  $x^*$  is a strict local minimum. Conversely if  $x^*$  is a critical point such that  $Hf(x^*)$  is positive definite the left hand side of (4.9) remains strictly positive for all  $t \neq 0$  sufficiently small and all  $v \neq 0 \in T_{x^*}$  which shows that  $x^*$  is a strict local minimum. ■

Remark 4.1 (Generic properties). Morse functions form an open dense subset of the space of  $\mathcal{C}^\sigma$  real-valued functions on  $C$  for  $2 \leq \sigma \leq +\infty$  (see Hirsch (Ref. 14)) for a precise definition of the topology). In other words, by a "slight perturbation" it is always possible to ensure that  $f$  is a Morse function. Such a property is said to be generic (See Golubitsky, Guillemin (Ref. 15)).

Incidentally the other assumption we made, namely 0 is a regular value of the  $\mathcal{C}^\sigma$  map  $C : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is also generic for  $1 \leq \sigma \leq +\infty$ . In terms of differential topology, the assumption amounts to said that  $c$  is transverse to 0 on the manifold  $C$  ( $c \nparallel \{0\}$ ). The genericity of the assumption follows from the transversality theorem (See e.g. (Ref. 14)).

These remarks show that our assumptions hold for "most"  $\mathcal{C}^\sigma$  problems (1.1). ■

Remark 4.2 The proofs of theorems 4.1 and 4.2 use explicitly the Riemannian structure of  $C$  since they introduce geodesics. They generalize similar results for the unconstrained case which actually make use of the Euclidean structure of  $\mathbb{R}^n$ . Such results hold however on a manifold without a Riemannian structure and can be established using a local coordinate system and showing they are independent of its choice. ■

Remark 4.3 The reader may at this point wonder how these results relate to the classical formulation of the optimality conditions for problem (1.1) in terms of Lagrange multipliers. The first order optimality condition (4.3) is equivalent to

$$(4.9) \quad f'(x^*)(v) = \langle \nabla f(x^*), v \rangle_n = 0 \quad \forall v \in T_{x^*}$$

by definition of the gradient  $\nabla f$  of  $f$  in the Euclidean space  $\mathbb{R}^n$ ; condition (4.9) expresses that  $\nabla f(x^*)$  is orthogonal to  $T_{x^*}$  hence belongs to  $\mathcal{N}(A_{x^*}^T)$ . Thus there exists  $\lambda^* \in \mathbb{R}^m$  such that

$$(4.10) \quad \nabla f(x^*) + \nabla c(x^*) \cdot \lambda^* = 0,$$

the classical form for the first-order optimality condition in term of the Lagrangian

$$(4.11) \quad \ell(x, \lambda) = f(x) + \langle \lambda, c(x) \rangle_m$$

Given  $x \in C$ , let

$$(4.12) \quad \lambda(x) = - (\nabla c(x)^T \cdot \nabla c(x))^{-1} \nabla c(x)^T \nabla f(x);$$

clearly  $\nabla_x \ell(x, \lambda(x)) \in T_x$  and  $\lambda^* = \lambda(x^*)$  satisfies (4.10). At the critical point  $x^*$

$$(4.13) \quad Hf(x^*)(v, v) = \nabla_{xx}^2 \ell(x^*, \lambda^*)(v, v) \quad \forall v \in T_{x^*}$$

and the second-order optimality condition of theorem 4.2 can be stated in term of the positive definiteness of the restriction to  $T_{x^*}$  of the Hessian of the Lagrangian  $\nabla_{xx}^2 \ell(x^*, \lambda^*)$ . (See Luenberger (Ref. 13)). ■



#### 4.2. - Gradient of a function on C

Let  $\gamma$  the Riemannian structure on C. We define the gradient on C of f at x as the tangent vector  $\nabla_C^\gamma f(x) \in T_x$  such that

$$(4.14) \quad \gamma_x (\nabla_C^\gamma f(x), v) = D_C f(x)(v) \quad \text{for all } v \in T_x.$$

The vector  $\nabla_C^\gamma f(x)$  clearly depends upon the Riemannian metric. Notice that  $\nabla_C^\gamma f$  is a vector field on C called the gradient field.

In the local coordinate system around x

$$(4.15) \quad z(y) = Z_x(y - x) \quad \text{for all } y \in U_x \cap C,$$

there are two natural Riemannian metrics as discussed in section 3.1. If we use the metric  $\gamma_x^Z$  defined in (3.5) as the scalar product in  $\mathbb{R}^{n-m}$  we obtain

$$(4.16) \quad \nabla_C^Z f(x) = Z_x^{-1} g_x^Z \quad \text{with } g_x^Z = Z_x^{-T} \nabla f(x) \in \mathbb{R}^{n-m}$$

while if we use the metric  $\gamma_x^E$  induced by the scalar product in  $\mathbb{R}^n$  defined in (3.4) we obtain

$$(4.17) \quad \nabla_C^E f(x) = Z_x^{-1} g_x^E \quad \text{with } g_x^E = (Z_x^{-T} Z_x^{-1})^{-1} g_x^Z.$$

We call  $g_x^Z$  and  $g_x^E$  respectively the reduced gradient and the Euclidean reduced gradient; they coincide obviously iff  $Z_x^{-T} Z_x^{-1} = I$ , which gives an argument in favor of the coordinate system presented in Example 2.2. The terminology reduced gradient introduced by Abadie (Ref. 5) expresses the reduction of dimension achieved by considering vectors in  $\mathbb{R}^{n-m}$  and is employed here in a more general sense than in Gabay-Luenberger (Ref. 1) where only the partitioned coordinate system of Example 2.1 was considered. Notice that

$$(4.18) \quad \nabla_C^E f(x) = P_x^E \nabla f(x),$$

where  $P_x^E$  is the orthogonal projector onto  $T_x$  in the Euclidean space  $\mathbb{R}^n$ ;  $\nabla_C^E f(x)$  will be called the (Euclidean) projected gradient following Rosen (Ref. 3). The vector  $\nabla_C^Z f(x)$  can also be interpreted as the orthogonal projection of  $\nabla f$  onto  $T_x$  but with respect to the scalar product

$$(4.19) \quad \langle\langle v, w \rangle\rangle = \langle S_x v, S_x w \rangle_n,$$

with  $S_x$  defined by (2.5).

Remark 4.4 Another argument in favor of the local coordinate system defined in Example 2.2 arise from the relationship between the projected gradient  $\nabla_C^E f(x)$  and the gradient  $\nabla_x \ell(x, \lambda(x))$  of the Lagrangian. Both vectors are in  $T_x$  provided  $\lambda(x)$  is given by (4.12) which can be expressed (and computed) in term of the pseudo-inverse (2.19) of  $A_x = \nabla c(x)^T$ ,

$$(4.20) \quad \lambda(x) = - A_x^{+T} \nabla f(x).$$

We then have

$$(4.21) \quad \nabla_C^E f(x) = P_x^E \nabla f(x) = \nabla_x \ell(x, \lambda(x)),$$

formula used by Luenberger (Ref. 2 ) to define the gradient of  $f$  under constraints (see also Hestenes (Ref. 16)).

If we use a general right inverse  $A_x^-$  and define

$$(4.22) \quad \mu(x) = - A_x^{-T} \nabla f(x),$$

then by (2.10)

$$(4.23) \quad \nabla_x \ell(x, \mu(x)) = Z_x^T Z_x^- \nabla f(x) = Z_x^T g_x^Z$$

which differs from  $\nabla_C^Z f(x)$  and actually does not belong to  $T_x$  except if  $Z_x^- = Z_x^T$ , i.e. if the columns of  $Z_x^T$  are orthonormal. ■

#### 4.3. - Steepest Descent Method along geodesics

To find a local minimum on  $\mathbb{R}^n$  of a continuously differentiable function  $f$ , the steepest descent method generates, starting from an initial estimate  $x^0$ , successive approximations according to the iteration

$$(4.24) \quad x^{k+1} = x^k + t_k p^k, \quad k = 0, 1, \dots,$$

where  $p^k$  is the direction which minimizes  $\langle \nabla f(x^k), p \rangle_n / \|p\|_n$ , namely

$$(4.25) \quad p^k = - \nabla f(x^k),$$

called the direction of steepest descent, and the stepsize  $t_k$  is a positive scalar selected for instance as the first local minimum on  $\mathbb{R}^+$  of the function

$j(t) = f(x^k + t p^k)$ ; we say that  $t_k$  is determined by a exact line search and denote the solution by

$$(4.26) \quad t_k = \operatorname{Argmin} \{f(x^k + t p^k) \mid t \geq 0\}.$$

This method can be generalized to find a local minimum of  $f$  on a Riemannian manifold  $C$ . The direction of steepest descent for  $f$  on  $C$  at  $x^k$  is given by

$$(4.27) \quad p^k = - \nabla_C^\gamma f(x^k),$$

which minimizes  $\gamma_{x^k}(\nabla_C^\gamma f(x^k), p) / \|p\|_\gamma$  for all  $p \in T_{x^k}$  ( $\|p\|_\gamma = (\gamma_{x^k}(p, p))^{1/2}$ ).

Proposition 3.3 shows that the geodesics of  $C$  play the role of the straight lines in  $\mathbb{R}^n$ ; we thus define the steepest descent method along geodesics as the iteration

$$(4.28) \quad x^{k+1} = \exp_{x^k}(t_k p^k),$$

where  $p^k$  is defined by (4.27) and the stepsize  $t_k$  is determined by an exact geodesic search,

$$(4.29) \quad t_k = \operatorname{Argmin} \{f(\exp_{x^k}(t p^k)) \mid t \in \mathbb{R}^+\}.$$

Algorithm (4.27) (4.28) (4.29) has been first introduced and analyzed by Luenberger (Ref. 2) who explicitly used the Riemannian structure  $\gamma^E$  on  $C$  induced by the Euclidean structure of  $\mathbb{R}^n$ . As noticed in (4.18)  $\nabla_C^E f(x^k)$  is then the orthogonal projection the gradient  $\nabla f(x^k)$  on  $T_{x^k}$ ; hence the terminology : gradient projection method along geodesics. In his paper, Luenberger established the global convergence of the algorithm to a critical point of  $f$  on  $C$  and estimated the speed of convergence in the neighborhood of a critical point which is a strict local minimum. Lichnewski (Ref. 4) has recently proposed the algorithm (4.27) (4.28) (4.29) for a general Riemannian manifold and established similar results; he also studied a conjugate-gradient version of it.

We give a global convergence theorem which generalizes the classical results for the method in  $\mathbb{R}^n$  (see Polak (Ref. 17), Ortega-Rheinboldt (Ref. 18)). We first need to introduce the following notation : let  $W_k$  denote the connected component containing  $x^k$  of the level set  $\{x \in C \mid f(x) \leq f(x^k)\}$ .

**THEOREM 4.3** Assume that  $f$  is continuously differentiable and that  $W_0$  is compact.  
Then the sequence  $\{x^k\}$  constructed by the steepest descent method along  
geodesics (4.27) (4.28) (4.29) is well defined; it is either finite,  
terminating at a critical point, or is infinite and there is a subsequence  
converging to a critical point. If the critical values of  $f$  are distinct the  
whole sequence  $\{x^k\}$  converges to a critical point. ■

**Proof :** The compactness of  $W_0$  implies the compactness of the closed subsets  $W_k$   
 since the sequence  $\{f(x^k)\}$  is monotone and non increasing. If  $x^k$  is a critical  
 point,  $p^k = 0$  and the algorithm does not generate new approximations;  
 introducing a stopping test it can be terminated at iteration  $k$ . Assume now  
 that  $x^k$  is not a critical point : hence

$$(4.30) \quad \gamma_{x^k}(\nabla_C^\gamma f(x^k), p^k) = -\gamma_{x^k}(\nabla_C^\gamma f(x^k), \nabla_C^\gamma f(x^k)) < 0.$$

Denote the arc of geodesic starting from  $x^k$  with tangent  $p^k$  by

$$(4.31) \quad x(t) = \exp_{x^k}(t p^k),$$

and introduce the function  $j : \mathbb{R}^+ \rightarrow \mathbb{R}$  defined by

$$(4.32) \quad j(t) = f[x(t)].$$

Let  $\bar{t} = \limsup J$  where the set  $J$  is defined by

$$J = \{t > 0 \mid x(t) \text{ is defined and } j(t) < j(0) = f(x^k)\};$$

by (4.30) the set  $J$  is not empty and since  $W_k$  is compact, Proposition 3.3  
 implies that either  $\bar{t} = +\infty$  and  $x(t) \in W_k$  for all  $t \in [0, +\infty)$  or  
 $\bar{t}$  is finite,  $f[x(\bar{t})] = f(x^k)$  and  $x(t) \in W_k$  for all  $t \in [0, \bar{t}]$ .

In both cases the stepsize rule (4.29) is well defined :  $t_k \in (0, \bar{t})$ .

Observe that  $t_k$  satisfies

$$(4.34) \quad j'(t_k) = \gamma_{x^{k+1}}(\nabla_C^\gamma f(x^{k+1}), \pi_0^{t_k}(p^k)) = \gamma_{x^k}(\pi_{t_k}^0(\nabla_C^\gamma f(x^{k+1})), p^k) = 0$$

where  $\pi_0^{t_k}(p^k)$  is the vector of  $T_{x^{k+1}}$  obtained from  $p^k$  by parallel translation  
 along the curve  $x(\cdot)$  from  $x^k$  to  $x^{k+1}$  (see proposition 3.1); the second  
 equality of (4.34) results from the definition of the geodesic  $x(\cdot)$ .

Observe also that given  $\alpha \in (0, \frac{1}{2})$  the equation

$$(4.35) \quad \gamma_{x^k}(\pi_t^0(\nabla_C^\gamma f(x(t))), p^k) = \alpha \gamma_{x^k}(\nabla_C^\gamma f(x^k), p^k)$$

has a smallest solution  $\hat{t} \in (0, t_k)$  and

$$\gamma_{x^k}(\pi_t^0(\nabla_C^\gamma f(x(t))), p^k) < \alpha \gamma_{x^k}(\nabla_C^\gamma f(x^k), p^k) \quad \text{for all } t \in [0, \hat{t}).$$

Applying the mean value theorem we obtain the estimate

$$f(x^{k+1}) - f(x^k) < f(x(\hat{t})) - f(x^k) < -\alpha \hat{t} \|\nabla_C^\gamma f(x^k)\|_\gamma^2.$$

The sequence  $\{f(x^k)\}$  is monotone decreasing and is bounded from below since the continuous function  $f$  attains its minimum on the compact  $W_0$ ; hence it converges to a limit. Let  $\{x^{k_i}\}$  a subsequence of  $\{x^k\}$  converging to  $x^* \in W_0$ . Suppose that  $x^*$  is not a critical point, i.e.  $\|\nabla_C^\gamma f(x^*)\|_{\gamma_{x^*}} = \delta > 0$ ;

The continuity of the Riemannian metric  $\gamma$  and of the gradient field  $\nabla_C^\gamma f$  implies that  $\|\nabla_C^\gamma f(x^{k_i})\|_\gamma \geq \delta/2$  for all  $i > I$ ; hence

$$f(x^{k_{i+1}}) < f(x^{k_i+1}) < f(x^{k_i}) - \alpha \hat{t} \delta^2/4,$$

which contradicts the fact that  $\{f(x^{k_i})\}$  converges to  $f(x^*)$ . Thus  $x^*$  is a critical point.

Finally suppose that  $x^*$  and  $x^{**}$  are distinct accumulation points of the sequence  $\{x^k\}$  in  $W_0$ :  $x^*$  and  $x^{**}$  are critical points of  $f$ . Since  $\{f(x^k)\}$  converges we must have  $f(x^*) = f(x^{**})$ , which is impossible if the critical values of  $f$  are distinct. ■

Remark 4.5. The proof shows that theorem 4.3 holds if we replace the stepsize rule (4.29) by rule (4.34) or rule (4.35) or the following rule: find  $t_k = 2^{-\ell} \cdot t$  where  $t$  is an initial guess and  $\ell$  is the smallest integer satisfying for a given  $\alpha \in (0, \frac{1}{2})$ ,

$$(4.36) \quad f(x(t_k)) \leq f(x^k) - \alpha t_k \|\nabla_C^\gamma f(x^k)\|_\gamma^2.$$

By analogy with stepsize selection rules for unconstrained minimization we call (4.34), (4.35) and (4.36) respectively the Curry, Altman and Armijo principles (see Ortega-Rheinhold (Ref. 18)). ■

Remark 4.6. A convergence result similar to theorem 4.3 is established by Lichnewski (Ref. 4) for a  $\mathcal{C}^2$  function  $f$ . We prefer to establish the result assuming only  $\mathcal{C}^1$  differentiability to extend the classical convergence theory of steepest descent methods to the constrained situation. In his paper Lichnewski introduces a special procedure to handle the case where the algorithm reaches a neighborhood of a critical point which is a saddle-point and not a local minimum; he generates in this case a new direction of descent which can be called a direction of negative curvature using the terminology recently introduced by Mc Cormick (Ref. 19), Moré-Sorensen (Ref. 20) for the unconstrained case. It is then possible to show that the modified algorithm generates a sequence converging always to a local minimum if  $f$  is a Morse function and has distinct critical values, which is true for "most"  $\mathcal{C}^2$  functions (see Remark 4.1). ■

We finally give an estimate of the speed of convergence.

THEOREM 4.4. Assume that  $f$  is a  $\mathcal{C}^\sigma$  Morse function on the Riemannian manifold  $C$  with  $\sigma \geq 3$ . Suppose that the steepest descent method along geodesics (4.27) (4.28) (4.29) generates a sequence  $\{x^k\}$  converging to a critical point  $x^*$  in  $C$  such that the Hessian form  $Hf(x^*)$  satisfies

(4.37)  $m \|v\|_{\gamma^*}^2 \leq Hf(x^*)(\dot{v}, v) \leq M \|v\|_{\gamma^*}^2$  for all  $v \in T_{x^*}$ , where  $m$  and  $M$  are two positive scalars. Then the sequence  $\{x^k\}$  is linearly convergent and

$$(4.38) \quad \lim_{k \rightarrow \infty} \frac{f(x^{k+1}) - f(x^*)}{f(x^k) - f(x^*)} \leq \left( \frac{M - m}{M + m} \right)^2. \quad \blacksquare$$

See Lichnewski (Ref. 4) for a proof (or Luenberger (Ref. 2) for the special case of the Riemannian metric induced by the Euclidean structure of  $\mathbb{R}^n$ ). This result extends the estimate of rate of convergence for the steepest descent method in  $\mathbb{R}^n$  (See Luenberger (Ref. 13)). See also Gabay-Luenberger (Ref. 1) where this estimate was obtained for the special normal coordinate system described in Example 2.1.

Notice that the steepest descent method depends upon the Riemannian metric on  $C$  as well as the scalars  $m, M$  defined by (4.57) and giving the estimate of the speed of convergence.

#### 4.4. - Newton's Method along geodesics

Let  $x^*$  be a local unconstrained minimum of a  $\mathcal{C}^\sigma$  function  $f$  ( $\sigma \geq 3$ ) such that the Hessian form  $\nabla^2 f(x^*)$  is positive definite and let  $U_*$  be a neighborhood of  $x^*$  such that, for all  $x \in U_*$ ,  $\nabla^2 f(x)$  is positive definite. Starting from  $x^0 \in U_*$ , Newton's method for unconstrained minimization generates a sequence of approximations of  $x^*$  according to the iteration

$$x^{k+1} = x^k - (F_k)^{-1} \nabla f(x^k),$$

where  $F_k$  denotes the (non singular) symmetric matrix defining the Hessian form  $\nabla^2 f(x^k)$ ; provided  $U_*$  is sufficiently small the iterates remain in  $U_*$  and the sequence  $\{x^k\}$  is well-defined and converges to  $x^*$  with a quadratic rate of convergence (see e.g. Ortega-Rheinboldt (Ref. 18)).

The extension of Newton's method to the minimization on a manifold  $C$  presents a major difficulty since it is not possible to define the Hessian form of  $f$  on  $C$  outside of a critical point (see Section 4.1). We can however define at a non-critical point  $x \in C$  a quadratic form on the tangent space  $T_x$  by exploiting the Riemannian structure of the manifold  $C$ .

Let  $v \in T_x$  and consider the geodesic curve  $x(\cdot)$  starting from  $x$  and tangent to  $v$ ,

$$(4.39) \quad x(t) = \exp_x(tv).$$

We proceed like for the definition of the Hessian but let now  $V$  be the parallel vector field along the curve  $x(\cdot)$  such that  $V(0) = v$ ; it is unique by Proposition 3.1. Define the  $\mathcal{C}^\sigma$  function  $j : \mathbb{R}^+ \rightarrow \mathbb{R}$  by

$$(4.40) \quad j(t) = f(x(t)).$$

By definition of the gradient of  $f$  on  $C$  with respect to the Riemannian metric  $\gamma$  we have

$$(4.41) \quad j'(t) = \gamma_{x(t)}(\nabla_C f(x(t)), V(t))$$

and

$$j''(t) = \gamma_{x(t)}\left(\frac{D}{dt} \nabla_C f(x(t)), V(t)\right) + \gamma_{x(t)}(\nabla_C f(x(t)), \frac{DV}{dt}(t))$$

$$(4.42) \quad = \gamma_{x(t)} \left( \frac{D}{dt} \nabla_C f (x(t)), \dot{V}(t) \right),$$

since  $\dot{V}$  is a parallel vector field. We define  $F(x) : T_x \times T_x \rightarrow \mathbb{R}$  by

$$(4.43) \quad F(x) (v, v) = j''(0) = \gamma_x (D_v \nabla_C f, v) \text{ for all } v \in T_x.$$

The covariant derivative  $D_v \nabla_C f$  of the gradient field  $\nabla_C f$  along the vector  $v$  is linear in  $v$  (see (3.8a)); hence  $F(x)$  is a quadratic form. The regularity of the solution (4.39) of the differential equations (3.15) defining the geodesic with respect to the initial conditions implies that  $F$  is  $\mathcal{C}^{\sigma-2}$  differentiable with respect to  $x$ . Notice that at a critical point  $x^*$ ,  $F(x^*)$  coincides with the Hessian form defined by (4.8) :

$$(4.44) \quad F(x^*) (v, v) = Hf(x^*) (v, v) \quad \text{for all } v \in T_{x^*}.$$

Define the map  $F_x : T_x \rightarrow T_x$  by

$$(4.45) \quad F(x) (v, v) = \gamma_x (F_x v, v) \quad \text{for all } v \in T_x;$$

notice that  $F_x$  is simply the Riemannian connexion of the gradient field at  $x$ ,

$$(4.46) \quad F_x = \tau_x (\nabla_C f).$$

It is a linear isomorphism of  $T_x$  and is self adjoint (with respect to the Riemannian metric  $\gamma_x$ ); denote by  $(F_x)^{-1}$  its inverse (defined on  $T_x$ ).

Let  $x^*$  be a non degenerated local minimum of  $f$  on  $C$ ; by theorem 4.2 the Hessian  $Hf(x^*)$  is positive definite. Identity (4.44) together with the continuity of  $F$  show that there exists a neighborhood  $U_* \cap C$  of  $x^*$  in  $C$  such that  $F(x)$  is positive definite for all  $x \in U_* \cap C$ . Starting from  $x^0 \in U_* \cap C$  Newton's method along geodesics generates the sequence of successive approximations  $\{x^k\}$  according to the iteration.

$$(4.47) \quad x^{k+1} = \exp_{x^k} (p^k),$$

where the "Newton's direction"  $p^k \in T_{x^k}$  is given by

$$(4.48) \quad p^k = - (F_{x^k})^{-1} \nabla_C f(x^k).$$

The method is well-defined provided  $U_* \cap C$  is chosen small enough so that the



sequence  $\{x^k\}$  remains in this neighborhood, and  $\{x^k\}$  converges to  $x^*$  quadratically.

**THEOREM 4.5.** Assume that  $f$  is a  $\mathcal{C}^\sigma$  Morse function on the Riemannian manifold  $C$  with  $\sigma \geq 3$ . Let  $x^*$  be a local minimum of  $f$  on  $C$ . Then there exists a neighborhood  $U_* \cap C$  of  $x^*$  in  $C$  such that if  $x^0 \in U_* \cap C$  the Newton's method along geodesics generates a sequence  $\{x^k\}$  in  $U_* \cap C$  converging to  $x^*$  and there exists a constant  $K$  such that

$$(4.49) \quad \delta(x^{k+1}, x^*) \leq K(\delta(x^k, x^*))^2 \quad \text{for all } k = 0, 1, \dots,$$

where  $\delta(.,.)$  stands for the Riemannian distance. ■

Proof : Given  $x \in C$  and  $v \in T_x$  define the  $\mathcal{C}^{\sigma-1}$  function  $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$  by

$$(4.50) \quad \phi(t) = \gamma_{x(t)}(\nabla_C f(x(t)), W(t)),$$

where  $W(t)$  is a parallel vector field along the curve  $x(.)$  given by (4.39). Notice that

$$(4.51) \quad \phi'(t) = \gamma_{x(t)}(F_{x(t)} V(t), W(t)),$$

where  $V(t)$  is the parallel vector field along  $x(.)$  such that  $V(0) = v$ .

Consider first the geodesic curve  $x(.)$  starting at  $x^k$  and tangent to  $p^k$  given by (4.48). Taylor's formula for the function  $\phi(t)$ ,  $\phi(1) = \phi(0) + \phi'(0) + \int_0^1 [\phi'(t) - \phi'(0)] dt$ , yields

$$(4.52) \quad \gamma_{x^{k+1}}(\nabla_C f(x^{k+1}), W(1)) = \int_0^1 \gamma_{x^{k+1}}(\pi_t^1 [F_x(t) - \\ - \pi_0^t F_{x^k} \pi_t^0] V(t), W(1)) dt;$$

since  $F_x$  is  $\mathcal{C}^{\sigma-2}$  differentiable on  $U_* \cap C$  we obtain the estimate

$$(4.53) \quad \|\nabla_C f(x^{k+1})\|_\gamma \leq \frac{L}{2} \|p^k\|_\gamma^2 \leq \frac{L}{2m^2} \|\nabla_C f(x^k)\|_\gamma^2,$$

if  $\|(F_x)^{-1}\| \leq 1/m$  for all  $x \in U_* \cap C$ .

Consider now the minimal geodesic joining  $x^{k+1}$  and  $x^*$ , which exists by Proposition 3.3; we can find  $q^{k+1} \in T_{x^{k+1}}$  such that the curve

$$x(t) = \exp_{x^{k+1}}(t q^{k+1})$$

satisfies  $x(0) = x^{k+1}$  and  $x(1) = x^*$ . For this choice of  $x(\cdot)$ , Taylor's formula for the function  $\varphi$  yields

$$(4.54) \quad \gamma_{x^*}(\nabla_C f(x^*), W(1)) = 0 = \gamma_{x^{k+1}}(\nabla_C f(x^{k+1}), W(0)) + \int_0^1 \gamma_{x(t)}(F_{x(t)} Q(t), W'(t)) dt;$$

assuming that  $m \|v\|_\gamma^2 \leq \gamma_x(F_x v, v) \leq M \|v\|_\gamma^2$  for all  $x \in U_* \cap C$

we obtain the estimate

$$(4.55) \quad m \|q^{k+1}\|_\gamma \leq \|\nabla_C f(x^{k+1})\|_\gamma \leq M \|q^{k+1}\|_\gamma.$$

Combining (4.53) and (4.55) for  $k$  and  $k-1$  and noticing that  $\delta(x^{k+1}, x^*) = \|q^{k+1}\|_\gamma$  we obtain

$$(4.56) \quad \delta(x^{k+1}, x^*) \leq \frac{L M}{2 m^3} (\delta(x^k, x^*))^2$$

which shows the quadratic convergence of  $\{x^k\}$  to  $x^*$  provided

$$\delta(x^0, x^*) \leq \frac{2 m^3}{L M}. \blacksquare$$

#### 4.5. - Quasi Newton Methods along Geodesics

It is possible to modify Newton's method along geodesics in order to have a globally convergent method, i.e. a method which constructs a sequence  $\{x^k\}$  converging to a critical point  $x^*$  from any initial approximation  $x^0 \in C$ . The modified method consists in generating the sequence  $\{x^k\}$  according to the iteration

$$(4.57) \quad x^{k+1} = \exp_{x^k}(t_k p^k),$$

where  $p^k$  is a descent direction defined by

$$(4.58) \quad p^k = - (G_k)^{-1} \nabla_C f(x^k),$$

and the stepsize  $t_k$  is chosen according to Armijo principle (4.36), i.e. given  $\alpha \in (0, \frac{1}{2})$  let  $t_k = 2^{-\ell}$  with  $\ell$  the smallest integer such that

$$(4.59) \quad f(\exp_x^k(2^{-l} p^k)) \leq f(x^k) - \alpha 2^{-l} \gamma_{x^k} (\nabla_C f(x^k), (G_k)^{-1} \nabla_C f(x^k)).$$

The method (4.57) (4.58) (4.59) generalizes to the constrained case the variable metric methods in  $\mathbb{R}^n$ .

We now specify how the operator  $G_k$  of  $T_k$  must be chosen in order that  $p^k$  approximates the Newton direction (4.48). Observe that (4.52) yields for Newton's method the equation

$$(4.60) \quad \nabla_C f(x^{k+1}) = \pi_0^1(\nabla_C f(x^k)) + \int_0^1 \pi_t^1 F_{x(t)} \pi_0^t(p^k) dt.$$

Define the linear operator  $\bar{F}_{k+1} : T_{k+1} \rightarrow T_{k+1}$  by

$$(4.61) \quad \bar{F}_{k+1} = \int_0^1 \pi_t^1 F_{x(t)} \pi_1^t dt$$

Equation (4.60) can be rewritten as

$$(4.62) \quad \nabla_C f(x^{k+1}) - \pi_0^1(\nabla_C f(x^k)) = \bar{F}_{k+1} \pi_0^1(p^k);$$

it is thus natural in order to approach Newton's method to require that the operator  $G_{k+1}$  satisfies the Quasi-Newton equation for constrained minimization

$$(4.63) \quad \nabla_C f(x^{k+1}) - \pi_0^{t_k} \nabla_C f(x^k) = G_{k+1} \pi_0^{t_k}(t_k p^k),$$

where we have slightly modified (4.62) to take into account the presence of the stepsize  $t_k$  (we recall that  $\pi_0^{t_k}$  denote the parallel translation along the curve  $x(t) = \exp_k^x(tp^k)$  from  $x(0) = x^k$  to  $x(t_k) = x^{k+1}$ ). Many updating formulae for  $G_k$  can be proposed which satisfy the Quasi-Newton equation and generalize to the constrained case the constellation of Quasi-Newton methods for the minimization of  $f$  on  $\mathbb{R}^n$  (See the excellent survey by Dennis-Moré (Ref. 21)). We can specialize further the choice of the updating scheme by requiring some additional properties. For instance we want that  $p^k$  be always a descent direction in order to find an admissible stepsize  $t_k$  satisfying (4.59); hence the updating formula should generate from a positive definite operator  $G_k$  on  $T_k$  a new operator  $G_{k+1}$  which is positive definite on  $T_{k+1}$ . In order to approximate properly the self-adjoint operator  $(F_k)^{-1}$ , we can insist on having self adjoint operators  $(G_k)^{-1}$ ; hence the updating formula should preserve the symmetry of the operators  $(G_k)^{-1}$ . These two requirements in addition to the Quasi-Newton equation (4.63) are satisfied by a family of rank-two updating formula. Because of the recognized superiority of the Broyden-Fletcher-Goldfarb-Shanno updating scheme

for minimization in  $\mathbb{R}^n$ , we establish the corresponding formula for the constrained case. Denoting

$$(4.64) \quad s^k = \pi_o^{t_k}(t_k p^k) \in T_{k+1},$$

$$(4.65) \quad y^k = \nabla_C f(x^{k+1}) - \pi_o^{t_k} \nabla_C f(x^k) \in T_{k+1},$$

we define the operator  $G_{k+1} : T_{k+1} \rightarrow T_{k+1}$  by

$$(4.66) \quad G_{k+1} p = \tilde{G}_k p - \frac{\gamma_{k+1}(s^k, \tilde{G}_k p)}{\gamma_{k+1}(s^k, \tilde{G}_k s^k)} \tilde{G}_k s^k + \frac{\gamma_{k+1}(y^k, p)}{\gamma_{k+1}(y^k, s^k)} y^k$$

for all  $p \in T_{k+1}$ ,

with  $\tilde{G}_k = \pi_o^{t_k} G_k \pi_{t_k}^o$ . If  $G_k$  is positive definite on  $T_k$ ,  $G_{k+1}$  is positive definite on  $T_{k+1}$  iff

$$(4.67) \quad \gamma_{k+1}(y^k, s^k) > 0;$$

we call (4.66) the Generalized BFGS update formula for constrained minimization. Notice that the computation of  $G_{k+1}$  requires only first-order information,

namely the gradient at  $x^k$  and  $x^{k+1}$ , a definite advantage over the operator  $F_k$  used in Newton's method which involves second order information. To prove the global convergence of the method (4.57)(4.58)(4.59) with update (4.66) we must extend the very technical analysis of the BFGS method by Powell (Ref. 22) to the constrained case. This can be achieved provided there exists a constant  $M$  such that the inequality

$$(4.68) \quad \gamma_{k+1}(y^k, y^k) \leq M \gamma_{k+1}(y^k, s^k)$$

holds for all  $k$ ; we must also observe that the stepsize rule (4.59) implies

$$(4.69) \quad \gamma_{k+1}(\nabla_C f(x^{k+1}), \pi_o^{t_k} p^k) \geq \alpha' \gamma_k(\nabla_C f(x^k), p^k)$$

with  $0 < \alpha < \alpha' < 1$ . We can then establish a result similar to the one of Powell.

**THEOREM 4.6.** Assume that  $f$  is a  $\mathcal{C}^\sigma$  Morse function (with distinct critical values) on the Riemannian manifold  $C$  with  $\sigma \geq 2$  and that the level set  $W_0$  is compact. Let  $G_0$  be any positive definite self adjoint operator on  $T_0$ . Suppose that the sequence  $\{x^k\}$  constructed by the Quasi-Newton method along geodesics (4.57)(4.58)(4.59) with update (4.66) satisfies (4.68). Then the sequence  $\{x^k\}$  either terminates at or converges to a critical point. ■

The superiority of this Quasi-Newton method along geodesics over the steepest descent method along geodesics of section 4.3 is evidenced by the following result on its speed of convergence.

**THEOREM 4.7.** Assume that  $f$  is a  $\mathcal{C}^\sigma$  Morse function on the Riemannian manifold  $C$  with  $\sigma \geq 2$ . If the sequence  $\{x^k\}$  constructed by the Quasi-Newton method along geodesics converges to a critical point  $x^*$  such that the Hessian form  $Hf(x^*)$  is positive definite, then the sequence  $\{x^k\}$  is superlinearly convergent, i.e.

$$(4.69) \quad \lim_{k \rightarrow +\infty} \frac{\delta(x^{k+1}, x^*)}{\delta(x^k, x^*)} = 0. \quad \blacksquare$$

Proof : We simply outline the argument. Extending to the constrained case

Powell's technique (Ref. 22) we first show that  $\sum_{k=0}^{+\infty} \delta(x^k, x^*)$  is bounded.

We then use estimates of the type of Dennis-Moré (Ref. 23) to show that

$$(4.70) \quad \lim_{k \rightarrow +\infty} \frac{\| [G_k - F_k] p^k \|_\gamma}{\| p^k \|_\gamma} = 0.$$

This implies that after a finite number of iterations the stepsize

$$(4.71) \quad \tau_k = 1$$

satisfies the test (4.59) provided  $\alpha < \frac{1}{2}$  and the limit (4.69) holds. ■

**Remark 4.7.** Condition (4.70) expresses that the direction  $p^k$  given by (4.58) asymptotically approach the Newton direction (4.48). ■

**Remark 4.8.** We could similarly define a Generalized Davidon-Fletcher-Powell update; but the stepsize  $\tau_k$  of the corresponding Quasi-Newton method along geodesics must then be chosen by an exact geodesic search (4.29) in order to show the global and superlinear convergence. ■

## 5. - PRACTICAL IMPLEMENTATION

### 5.1. - Coordinate System

In practice the manifold  $C$  must be described by an atlas of local coordinate systems. Let  $x^k$  the current of any of the methods presented in Section 4 and let  $A_k$  be the  $m \times n$  Jacobian matrix of the map  $c$  defining  $C$ . In the neighborhood  $U_k \cap C$  of  $x^k$  in  $C$  we use the local coordinate system

$$(5.1) \quad z_k(x) = Z_k (x - x^k),$$

where  $Z_k$  is a  $(n-m) \times n$  matrix of right inverse  $Z_k^-$ , both chosen in term of  $A_k$  and its right inverse  $A_k^-$  as described in Proposition 2.1, namely

$$(5.2) \quad Z_k A_k^- = 0, \quad A_k Z_k^- = 0.$$

Let  $\theta_k : V \subset \mathbb{R}^{n-m} \rightarrow U_k \cap C$  denote the corresponding local parametrization of  $C$  around  $x^k$  defined in (2.4) such that

$$(5.3) \quad \theta_k(z_k(x)) = x \quad \text{for all } x \text{ in } C;$$

notice that  $\theta_k(0) = x^k$ . Theorem 2.2 shows that  $\theta_k$  is defined on the neighborhood  $V = B(0, 1/(2\beta\gamma\xi))$  of the origin in  $\mathbb{R}^{n-m}$ . In the local coordinate system (5.1) the tangent space  $T_k$  to  $C$  at  $x^k$  can be represented by  $\mathcal{R}(Z_k^-)$ ,

$$(5.4) \quad T_k = \{Z_k^- q \mid q \in \mathbb{R}^{n-m}\},$$

and a natural choice for the Riemannian metric at  $x^k$  is

$$(5.5) \quad \gamma_k(Z_k^- q, Z_k^- q') = \langle q, q' \rangle_{n-m} \quad \text{for all } q, q' \in \mathbb{R}^{n-m}$$

As noticed in Section 3.2, the coordinate system (5.1) is a normal coordinate system with respect to the Riemannian metric (5.5) and the geodesic curve  $x(\cdot)$  starting from  $x^k$  and tangent to  $v^k = Z_k^- q^k \in T_k$  is simply

$$(5.6) \quad x(t) = \theta_k(t q^k);$$

notice that if  $\|q^k\|_{n-m} = \alpha_k$ , the function  $x(t)$  is only defined for  $t \in [\tilde{t}, \tilde{t}]$  with

$$(5.7) \quad \tilde{\tau} = 1 / (2\alpha_k \beta \gamma \xi).$$

We can now describe the various descent methods of section 4 in the coordinate system (5.1). According to (4.16) the gradient of  $f$  on  $C$  at  $x^k$  can be expressed in term of the reduced gradient

$$(5.8) \quad g^k = Z_k^{-T} \nabla f(x^k).$$

The descent directions  $p^k$  defined by (4.25), (4.48), (4.58) can be represented by the general formula

$$(5.9) \quad p^k = - Z_k^{-1} H_k g^k,$$

where  $H_k$  is a symmetric matrix of dimension  $(n-m)$ ; the direction of steepest descent corresponds to  $H_k = I$ , the Newton's direction to  $H_k = Z_k (F_k)^{-1} Z_k^T$  and the Quasi-Newton direction to  $H_k = Z_k (G_k)^{-1} Z_k^T$ . According to (5.6) the typical iteration is

$$(5.10) \quad x^{k+1} = \theta_k (\tilde{\tau}_k (-H_k g^k))$$

where the stepsize  $\tilde{\tau}_k$  must now be selected on the interval  $[0, \tilde{\tau}]$ . Since  $\alpha_k = \|H_k g^k\|$  goes to 0, the upper bound  $\tilde{\tau}$  given by (5.7) increases and after a finite number of iterations  $\tilde{\tau}_k$  coincides with the stepsize  $\tau_k$  defined by any of the selection rules (4.29), (4.34), (4.35) or (4.36).

Remark 5.1 We can introduce a preconditioning constant matrix  $D$  of dimension  $(n-m)$  and define the Riemannian metric by  $\gamma_k(Z_k^{-1} q, Z_k^{-1} q') = \langle q, \frac{1}{2}(D + D^T)q' \rangle_{n-m}$  ■

Remark 5.2 In the Riemannian metric (5.5) the parallel translation

$\pi_o^{\tilde{\tau}_k} v^k$  of the vector  $v^k = Z_k^{-1} q^k \in T_k$  is the vector  $\hat{v}^k = Z_{k+1}^{-1} q^k \in T_{k+1}$ ; the Quasi-Newton equation (4.63) thus becomes in the coordinate system (5.1)

$$(5.11) \quad y^k = g^{k+1} - g^k = (H_{k+1})^{-1} s^k,$$

with  $s^k = -\tilde{\tau}_k H_k g^k$ , and the Generalized BFGS update formula (4.66) induces the rank-two update formula for the  $(n-m) \times (n-m)$  symmetric matrix  $H_k^{-1}$

$$(5.12) \quad H_{k+1}^{-1} = H_k^{-1} + \frac{y^k (y^k)^T}{\langle y^k, s^k \rangle} - \frac{H_k^{-1} s^k (s^k)^T H_k^{-1}}{\langle s^k, H_k^{-1} s^k \rangle}$$

Notice that the sequence of matrices  $H_k$  satisfies the formula

$$(5.13) \quad H_{k+1} = \left( I - \frac{s^k (y^k)^T}{\langle y^k, s^k \rangle} \right) H_k \left( I - \frac{y^k (s^k)^T}{\langle y^k, s^k \rangle} \right) + \frac{s^k (s^k)^T}{\langle y^k, s^k \rangle} \quad \blacksquare$$

## 5.2. - The Tangent-Restoration Approach

The constructive proof of Theorem 2.2 provides us with a practical scheme to compute  $x^{k+1}$  as

$$(5.14) \quad x^{k+1} = \theta_k (\tilde{t}_k p^k) = x^k - \tilde{t}_k Z_k^- H_k g^k + A_k^- w^k,$$

where  $w^k \in \mathbb{R}^m$  is chosen such that  $x^{k+1} \in C$ . Formula (5.14) can be interpreted geometrically in the following way :  $x^{k+1}$  is obtained by a tangent step from  $x^k$  to  $\tilde{x}^k = x^k - \tilde{t}_k Z_k^- H_k g^k$  in the affine tangent space  $x^k + T_k$  followed by a restoration step  $A_k^- w^k$  to enforce feasibility of  $x^{k+1}$ . Following (2.23), the restoration step is determined by taking  $w^k$  as the limit of the sequence  $w^{k,i}$  starting from  $w^{k,0} = 0$ , defined iteratively by

$$(5.15) \quad w^{k,i+1} = w^{k,i} - c(\tilde{x}^k + A_k^- w^{k,i}) \quad i = 0, 1, \dots;$$

since  $\tilde{t}_k \|p^k\| \leq 1/(2\beta\gamma\xi)$  the sequence  $\{w^{k,i}\}$  has a limit  $w^k$ .

The (approximate) geodesic search for the stepsize  $\tilde{t}_k$  can then be performed using a finite number of values for the parameter  $t$  and computing the restoration step by (5.15) from the finite number of points

$$(5.16) \quad \tilde{x}(t) = x^k - t Z_k^- H_k g^k$$

in the affine tangent space. The values for the parameter  $t$  can be obtained, starting from an initial guess in  $[0, \tilde{t}]$ , by an interpolation scheme (e.g. golden section, see Polak (Ref. 17)) to approximate the exact line search (4.29) or by successive halving to satisfy Armijo's rule (4.36).

For the steepest descent method ( $H_k = I$ ) it is shown in Gabay-Luenberger (Ref. 1) that an efficient initial guess can be obtained by the first local minimum on  $[0, \tilde{t}]$  of the Lagrangian function  $\ell(\tilde{x}(t), \mu^k)$  defined in (4.11) where  $\mu^k = -A_k^{-T} \nabla f(x^k)$  is the approximate Lagrange multiplier (4.22); the resulting approximation of the idealized steepest descent method converges linearly with



the optimal rate given by (4.38) in Theorem 4.4. A still simpler procedure is presented in Gabay (Ref. 24) and yields the same rate of convergence property.

For the Quasi-Newton method ( $H_k$  given by update formula (5.12)) an obvious initial guess is  $\text{Min}(1, \tilde{t})$ . By Theorem 4.7, after a finite number of iterations, the stepsize  $\tilde{t}_k = 1$  satisfies Armijo's rule (4.36) and the method converges superlinearly.

### 5.3. - Efficient Quasi-Newton methods for Constrained Minimization

We finally present the practical implementations of the Quasi-Newton method in the Tangent-Restoration approach for the specific coordinate systems envisaged in Examples 2.1 and 2.2.

(5.17) REDUCED QUASI-NEWTON METHOD. Given  $\epsilon_1, \epsilon_2 > 0$  (tolerance parameters),  $\tau > 0$  (estimate of  $1/2 \beta \gamma \xi$ ),  $\alpha \in (0, \frac{1}{2})$ ,  $x^k \in C$ ,  $H_k^{-1}$  a definite positive symmetric matrix of dimension  $(n-m)$  :

- i) Partition the Jacobian matrix  $A_k = [B, D]$ ; compute  $B^{-1}$ ;
- ii) Compute the reduced gradient  $g^k = \nabla_J f(x^k) - D^T (B^{-1})^T \nabla_I f(x^k)$ ;
- iii) If  $\|g^k\| \leq \epsilon_1$  then STOP, else  $\tilde{t} = \|H_k g^k\| \cdot \tau$ ; let  $t = \text{Min}(1, \tilde{t})$ ;
- iv) Tangent step :  $\tilde{x}(t) = x^k - t \begin{bmatrix} -B^{-1}D \\ I \end{bmatrix} H_k g^k$ ;  $w^{k,0} = 0$ ;
- v) Restoration step : while  $\|c(\tilde{x}(t) + \begin{bmatrix} B^{-1} \\ 0 \end{bmatrix} w^{k,i})\| > \epsilon_2$  do (5.15);  
 $x(t) = \tilde{x}(t) + \begin{bmatrix} B^{-1} \\ 0 \end{bmatrix} w^{k,i}$  ;
- vi) If  $(f(x(t)) \leq f(x^k) - \alpha t \langle g^k, H_k g^k \rangle)$  then  $x^{k+1} = x(t)$ , update  $H_k^{-1}$  according to (5.12); else assign  $\frac{t}{2}$  to  $t$  and go to step iv.

We prefer to update  $H_k^{-1}$  according to (5.12) rather than  $H_k$  directly according to (5.13) because the first scheme is numerically more "stable". To compute the tangent direction  $H_k g^k$  we then must solve a linear system; since the matrix of this system differs only by a rank-two correction from the one used in the previous iteration, the system can be solved in the order of  $(n-m)^2$  multiplications (see Gill-Murray (Ref. 25)).

A similar Reduced Quasi-Newton method was presented in Gabay-Luenberger (Ref. 1) using a Generalized Davidson-Fletcher-Powell update formula; it was however acknowledged that its performance was theoretically impaired by the approximate character of the linesearch in step vi. The Generalized BFGS update offers the superiority of not requiring exact linesearches (Theorem 4.7), which makes its choice particularly relevant for algorithm (5.17).

We now turn to the coordinate system defined in Example 2.2 which is a normal coordinate system with respect to the Riemannian metric on  $C$  induced by the Euclidean structure of  $\mathbb{R}^n$ . It allows to define a computationally efficient Quasi-Newton version of the gradient projection method.

(5.18) PROJECTED QUASI-NEWTON METHOD. Given  $\varepsilon_1, \varepsilon_2, \tau, \alpha, x^k \in C, H_k^{-1}$  defined as in (5.17) :

- i) Factorize the Jacobian matrix  $A_x = [L, 0]Q$  : partition  $Q = \begin{bmatrix} Q_1 \\ Q_2 \end{bmatrix}$  ;
- ii) Compute the Euclidean reduced gradient  $g^k = Q_2^T \nabla f(x^k)$  ;
- iii) If  $\|g^k\| \leq \varepsilon_1$  then STOP, else  $\tilde{t} = \|H_k g^k\| \cdot \tau$  ; let  $t = \text{Min}(1, \tilde{t})$  ;
- iv) Tangent step :  $\tilde{x}(t) = x^k - t Q_2^T H_k g^k$  ;  $w^{k,0} = 0$  ;
- v) Restoration step : while  $\|c(\tilde{x}(t) + Q_1^T L^{-1} w^{k,i})\| > \varepsilon_2$  do (5.15) ;  
 $x(t) = \tilde{x}(t) + Q_1^T L^{-1} w^{k,i}$  ;
- vi) If  $(f(x(t)) \leq f(x^k) - \alpha t \langle g^k, H_k g^k \rangle)$  then  $x^{k+1} = x(t)$ , update  $H_k^{-1}$  according to (5.12) ; else assign  $\frac{t}{2}$  to  $t$  and go to step iv.

This method generalizes to nonlinear equality constrained problems the Gill-Murray (Ref. 26) version of Goldfarb's method for linearly constrained problem (Ref. 27). Notice that it requires again the updating of the matrix  $H_k^{-1}$  of dimension  $(n-m)$  only. The difference between (5.17) and (5.18) lies mainly in their steps i; as noticed in section 2 the inversion by Gaussian elimination of the basic matrix  $B$  requires of the order of  $m^3/3$  multiplications while the QR decomposition of  $A_k$  can be obtained with approximately  $(n-\frac{m}{3}) m^2$  multiplications. However the conditioning of the matrices  $H_k$  generated by (5.18) is generally better than the one of the matrices generated by (5.17) (which may be severely affected by the chosen partition of  $A_k$ ).

Remark 5.3 : Methods (5.17) and (5.18) involve three scalar parameters  $\tau, \epsilon_1, \epsilon_2$ . The estimate  $\tau$  is introduced to guarantee that the restoration steps (v) provide a feasible point  $x(t)$ . In practice an exact value of  $\tau$  is unknown ; we generally start with a guessed approximation which is altered to a smaller value if the restoration phase fails in the process of iterations.

The idealized descent methods along geodesics correspond to the choice  $\epsilon_1 = \epsilon_2 = 0$ . They involve a (generally infinite) sequence of iterations requiring in their restoration step an infinite number of inner iterations. The positive tolerance parameters  $\epsilon_1, \epsilon_2$  are thus introduced to obtain implementable algorithms which, hopefully, terminate in a finite number of iterations and provide with an approximate solution. We must observe that the iterates  $x^k$  are then approximately feasible only (i.e.  $\|c(x^k)\| < \epsilon_2$ ), a situation which cannot be handled directly by our convergence theorems. Mukai and Polak (Ref. 30) have given a framework to establish convergence using such approximation schemes ; they essentially require that the parameters  $\epsilon_1, \epsilon_2$  be adapted carefully in the process of iterations to eventually reach an approximate solution. We describe an alternative approach in the second part of this paper where we propose a formally similar method defined for non-feasible points. ■

Remark 5.4 (A simple Arc Search) : In (Ref. 1) we proposed, for the implementation of the Reduced Newton's Method, to perform the approximation of the arc of geodesic

$$(5.19) \quad x(t) = \exp_{x^k}(tp^k)$$

using the second-order information available ; the sequential tangent restoration approach is then replaced by an approximate search along a parabola tangent to  $Z_k^{-k} p^k$  at  $x^k$ , followed by a restoration.

A similar device can be proposed in the framework of the Quasi-Newton method along geodesics, involving only first-order information. Suppose, for simplicity, that  $x^k$  is close enough to a solution so that  $t=1$  is a reasonable initial guess for the approximate search along the geodesic given by (5.19). Let  $p^k = -H_k g^k$  and  $q^k \in \mathbb{R}^m$ , the feasibility correction, be given by

$$(5.20) \quad q^k = -c(x^k + Z_k^{-k} p^k)$$

and consider the simple arc of curve

$$(5.21) \quad \tilde{x}(t) = x^k + t Z_k^{-k} p^k + \frac{t^2}{2} A_k^{-k} q^k \quad \text{for } t \in [0,1] ;$$

formula (5.21) defines an arc of parabola starting from  $x^k$  and tangent to  $Z_k^{-k} p^k$ .

The points  $\tilde{x}(t)$  for  $t > 0$  do not in general belong to  $C$  but can be viewed as an approximation of the geodesic obtained by performing only one step of the restoration phase (v).

We thus search the parabola (5.21) for the first  $t=2^{-\ell}$ ,  $\ell=0,1,\dots$  satisfying

$$(5.22) \quad f(\tilde{x}(t)) \leq f(x^k) + \alpha t \langle g^k, p^k \rangle ;$$

a restoration is then performed like in (v) and the new iterate defined like in (vi). This approach presents the advantage of requiring less restoration phases within an iteration. ■

Remark 5.5. In two recent papers (Ref. 28, 29) Tanabe has proposed and analyzed a continuous version of respectively the projected and the reduced Quasi-Newton method. He however requires the updating of a full  $n \times n$  non symmetric matrix, which is performed by Broyden's rank-one formula. ■

## 6. - PROBLEMS WITH INEQUALITY CONSTRAINTS

We consider now nonlinear programs where appear nonlinear inequality constraints. For simplicity we assume that no equality constraints are present; the problem is then given as

$$(6.1) \quad \text{Min } \{f(x) \mid x \in \mathbb{R}^n \text{ s. t. } c_i(x) \leq 0 \quad i = 1, \dots, p\},$$

where  $f$  and  $c_i$  are  $\mathcal{C}^\sigma$  real-valued functions on  $\mathbb{R}^n$ . The constrained set

$$(6.2) \quad C = \{x \in \mathbb{R}^n \mid c_i(x) \leq 0 \quad i = 1, \dots, p\}$$

can be embedded in a submanifold of  $\mathbb{R}^{n+p}$  in the following way. Define for  $i = 1, \dots, p$  the  $\mathcal{C}^\sigma$  functions  $g_i : \mathbb{R}^{n+p} \rightarrow \mathbb{R}$  by

$$(6.3) \quad g_i(x, z) = c_i(x) + z_i^2 \quad \text{for all } x \in \mathbb{R}^n, z \in \mathbb{R}^p,$$

where  $z_i$  is the  $i$ -th component of  $z$ , and consider the set  $G = g^{-1}(0) \subset \mathbb{R}^{n+p}$ . Given  $x \in C$  the index set of active constraints at  $x$  is the subset of  $P = \{1, \dots, p\}$  defined by

$$(6.4) \quad I(x) = \{i \in P \mid c_i(x) = 0\};$$

thus the corresponding point  $(x, z) \in G$  is such that  $z_i = 0$  for  $i \in I(x)$ , while  $z_i \neq 0$  for  $i \in P \setminus I(x)$ . Let  $m(x)$  denote the cardinality of  $I(x)$ . The Jacobian matrix of the map  $g$  is a  $p \times (n+p)$  matrix  $J_{x,z}$  partitioned as

$$(6.5) \quad J_{x,z} = [A_x, \Delta_z],$$

where  $A_x$  is the  $p \times n$  Jacobian matrix of the map  $c$  at  $x$  and  $\Delta_z$  is a  $p \times p$  diagonal matrix of diagonal entries  $2 z_i$ . Assume that for all  $x \in C$  the gradients of the  $m(x)$  active constraints  $\nabla c_i(x)$  for  $i \in I(x)$  are linearly independent. Then the matrix  $J_{x,z}$  is of full rank  $p$  and by Theorem 2.1 the set  $G$  is a  $\mathcal{C}^\sigma$  differential submanifold of  $\mathbb{R}^{n+p}$  of dimension  $n$ .

Problem (6.1) is thus equivalent to

$$(6.6) \quad \begin{array}{ll} \text{Min} & f(x), \\ (x,z) \in G & \end{array}$$

minimization of the  $\mathcal{C}^\sigma$  differentiable function  $f$  over the differential manifold  $G$  of class  $\sigma$ . Notice however that the problem is now formulated in the enlarged space  $\mathbb{R}^{n+p}$ . A naive approach, known in the mathematical programming literature as the active constraints strategy, consists in considering at each point  $x \in C$  only the active constraints and implementing the previous descent methods in a local coordinate system around  $x$  of the (variable) submanifold

$$(6.7) \quad C(x) = \{x \in \mathbb{R}^n \mid c_i(x) = 0 \quad i \in I(x)\}.$$

It turns out however that such a coordinate system cannot be used as a (partial) coordinate system for  $G$  around the corresponding point  $(x, z)$ ; in fact there exist in general points  $y$  in a neighborhood of  $x$  in  $C(x)$  such that  $c_i(y) > 0$  for some  $i$ , i.e.  $y \notin C$ . In order to apply efficiently the descent methods presented in this paper to problem (6.1) through formulation (6.6) we must design convenient local coordinate systems of  $G$  which exploit the separable structure of the functions  $g_i$ .

Another approach would consist in adapting directly our methods to the manifold with boundary  $C$  of  $\mathbb{R}^n$  by considering for instance local coordinate systems mapping a neighborhood  $U_x \cap C$  onto a neighborhood of the origin in a finite dimensional halfspace (See Hirsch (Ref. 14 § 1.4)).

## 7. CONCLUSIONS.

In this paper we have developed two distinct themes, a theoretic set-up using the geometry of manifolds and a computation-oriented analysis, and shown how they could be usefully interrelated. We have presented a geometric framework for studying nonlinear programming problems which has enabled us to generalize in an intrinsic manner the analysis and methods of unconstrained minimization to the constrained case. We have then specified the degrees of freedom of this general set-up to our advantage, using the theoretic approach as a guideline for the conception of efficient methods from the computational viewpoint.

We have in particular defined a family of descent methods along geodesics and presented their practical implementation. Such methods include most known primal methods for nonlinear programming and some new super-linear converging algorithms. They generate a sequence of feasible points, which requires at least conceptually, to solve a system of  $m$  nonlinear equations at each iteration. This inconvenience will be overcome in the second part of this paper.

## REFERENCES

- 1 D. GABAY and D.G. LUENBERGER "Efficiently Converging Minimization Methods Based on the Reduced Gradient", SIAM J. Control and Opt. 14, pp 42-61 (1976).
- 2 D.G. LUENBERGER, "The Gradient Projection Method along Geodesics", Management Science 18, pp 620-631 (1972).
- 3 J.B. ROSEN, "The Gradient Projection Method for Nonlinear Programming : Part II, Nonlinear Constraints", J. Soc. Indust. Appl. Math. 9, pp 514-522 (1961).
- 4 A. LICHNEWSKY, "Minimisation de fonctionnelles définies sur une variété par la méthode du gradient conjugué", Thèse de Doctorat d'Etat, Université Paris-Sud (1979).
- 5 J. ABADIE and J. GUIGOU, "Numerical Experiments with the GRG Method", Integer and Nonlinear Programming, J. Abadie ed., North-Holland, Amsterdam pp 529-536 (1970).
- 6 A. MIELE, H.Y. HUANG and J.C. HEIDEMAN, "Sequential Gradient-Restoration Algorithm for the Minimization of Constrained Functions", J.O.T.A. 4, pp 213-243 (1969).
- 7 J.W. MILNOR, Topology From the Differentiable Viewpoint, University Press of Virginia, Charlottesville (1965).
- 8 G.W. STEWART, Introduction to Matrix Computations, Academic Press, New York (1973).
- 9 V. GUILLEMIN and A. POLLACK, Differential Topology, Prentice Hall, Englewood Cliffs, N. J. (1974).
- 10 J.W. MILNOR, Morse Theory, Princeton University Press, Princeton N.J. (1969).
- 11 N.J. HICKS, Notes on Differential Geometry, Van Nostrand, Princeton, N.J. (1965).
- 12 R.L. BISHOP and R.J. CRITTENDEN, Geometry of Manifolds, Academic Press, New York (1964)

- 13 D.G. LUENBERGER, Introduction to Linear and Nonlinear Programming, Addison-Wesley, Reading, Mass (1973).
- 14 M.W. HIRSCH, Differential Topology, Springer, Heidelberg and New York (1976).
- 15 M. GOLUBITSKY and V. GUILLEMIN, Stable Mappings and their Singularities, Springer, Heidelberg and New York (1973).
- 16 M.R. HESTENES, Optimization Theory, the Finite Dimensional Case, Wiley New York (1975).
- 17 E. POLAK, Computational Methods in Optimization, A Unified Approach, Academic Press, New York (1971).
- 18 J.M. ORTEGA and N.C. RHEINBOLDT, Iterative Solution of Nonlinear Equations in Several Variables, Academic Press, New York (1970).
- 19 G.P. Mc CORMICK, "A Modification of Armijo's Step-size Rule for Negative Curvature", Math. Prog. 13, pp 111-115 (1977).
- 20 J.J. MORE and D.C. SORENSEN, "On the Use of Directions of Negative Curvature in a Modified Newton Method", Math. Progr. 16, pp 1-20 (1979).
- 21 J.E. DENNIS and J.J. MORE, "Quasi-Newton Methods, Motivation and Theory", SIAM Review 19, pp 46-89 (1977).
- 22 M.J.D. POWELL, "Some Global Convergence Properties of a Variable Metric Algorithm for Minimization Without Exact Line Searches", presented at the AMS/SIAM Symposium on Nonlinear Programming, New York (1976).
- 23 J.E. DENNIS and J.J. MORE, "A characterization of Superlinear Convergence and its application to Quasi-Newton Methods", Maths of Comp. 28, pp 549-560 (1974).
- 24 D. GABAY, "Efficient convergence of Implementable Gradient Algorithms and Stepsize Selection Procedures for Constrained Minimization", in International Computing Symposium 1975, E. Gelenbe and D. Potier eds., North-Holland, Amsterdam, pp 37-43 (1975).
- 25 P.E. GILL and W. MURRAY, "Quasi-Newton Methods for Unconstrained Optimization", J. Inst. Maths. Applics. 9, pp 91-108 (1972)



- 26 P.E. GILL and W. MURRAY, "Quasi-Newton Methods for Linearly Constrained Optimization", in Numerical Methods for Constrained Optimization, P.E. Gill and W. Murray eds., Academic Press, London (1974).
- 27 D. GOLDFARB, "Extension of Davidon's Variable Metric Method to Maximization Under Linear Inequality and Equality Constraints, SIAM J. Appl. Math. 17, pp 739-764 (1969).
- 28 K. TANABE, "A Geometric Method in Nonlinear Programming", STAN-CS-643; Stanford University (1977).
- 29 K. TANABE, "Differential Geometric Approach to Extended GRG Methods with Enforced Feasibility in Nonlinear Programming : Global Analysis", to appear in Recent Applications of Generalized Inverses, M.Z. Nashed ed..
- 30 H. MUKAI and E. POLAK, "On the Use of Approximations in Algorithms for Optimization Problems with Equality and Inequality Constraints", SIAM J. Numer. Anal., 13, pp. 674-693 (1978).

MINIMIZING A DIFFERENTIABLE FUNCTION  
OVER A DIFFERENTIAL MANIFOLD

PART II : QUASI-NEWTON METHODS WITH FEASIBILITY IMPROVEMENT \*

Daniel GABAY

Laboratoire d'Analyse Numérique. Université P. et M. Curie (Paris VI)

and

Institut National de Recherche en Informatique et Automatique

Domaine de Voluceau, 78150 LE CHESNAY (France).

ABSTRACT

We present a globally and superlinearly converging method for equality-constrained optimization requiring the updating of a reduced size matrix approximating a restriction of the Hessian of the Lagrangian. Each iterate is obtained by a search along a simple curve defined by a Quasi-Newton direction and a feasibility improving direction ; an exact penalty function is used to determine the stepsize. The method can be viewed as an efficient approximation to the Quasi-Newton along geodesics of (Ref. 1) where feasibility was enforced at each step. Its relation with multipliers methods and recursive quadratic programming methods is also investigated.

\* presented at the 10<sup>th</sup> International Symposium on Mathematical Programming  
Montreal, August 1979.

## I. INTRODUCTION

In a previous paper (Ref.1) we have studied a class of descent methods for the solution of the problem

$$(1.1) \quad \text{Min } \{f(x) \mid x \in \mathbb{R}^n \text{ s.t. } c_i(x) = 0 \text{ for } i=1, \dots, m\}$$

where the successive iterates  $x^k$  remain on the constraint manifold  $C = c^{-1}(0)$ . This class of methods includes most of the primal methods available today (e.g. gradient projection, reduced gradient methods). In particular, the framework adopted allows to define Quasi-Newton methods for the solution of (1.1) which require only the updating of an  $(n-m)$ -dimensional positive definite symmetric matrix approximating the Hessian of  $f$  on the manifold  $C$  at the minimum  $x^*$ ; such a method is shown to be superlinearly convergent. It requires however to search for an approximate local minimum of  $f$  along geodesic curves of the manifold  $C$ , endowed with a Riemannian metric. It is conceptually possible to achieve this scheme by two-phased iterations: a tangent step in the tangent space to  $C$  followed by a restoration phase to enforce the constraints. The resulting algorithm may be viewed as a sequence of solutions of  $m$  nonlinear equations, themselves performed iteratively (by a modified Newton's method); for large problems, like the optimal control of a nonlinear discrete time (implicit) dynamic system over a large number of periods, this second phase may require a much larger computational effort than the first phase.

In a practical implementation the restoration phase can only be performed approximately; the convergence theory becomes much more complex and requires the introduction of adaptive parameters for precision (see Mukai-Polak (Ref.2)).

It has been recognized since a long time that the constrained problem (1.1) could be solved by successive unconstrained minimizations of an augmented performance criteria. The penalty approach consists of successive minimizations of the functional

$$(1.2) \quad p(x, r_k) = f(x) + \frac{r_k}{2} ||c(x)||^2$$

for a monotonically increasing sequence of positive parameters  $r_k \rightarrow +\infty$ ; the resulting problems become increasingly ill-conditioned, which severely impairs the efficient minimization of (1.2). For this reason, the method of multipliers is preferred. At each iteration the augmented Lagrangian functional

$$(1.3) \quad \ell(x, \lambda^k, r_k) = f(x) + \langle \lambda^k, c(x) \rangle + \frac{r_k}{2} ||c(x)||^2,$$

obtained by adding a penalty term to the ordinary Lagrangian function

$$(1.4) \quad \ell(x, \lambda^k) = f(x) + \langle \lambda^k, c(x) \rangle,$$

is minimized and the parameters  $\lambda^k$  (and  $r_k$ ) updated. The resulting algorithm may be viewed as a dual method and, under some regularity assumptions (Ref.3), it may be shown to converge to a solution  $x^*$  without increasing  $r_k$  to  $+\infty$  (see the excellent survey by Bertsekas (Ref.4)).

Primal methods and multipliers methods require practically a sequence of (approximate) solutions of respectively a system of  $m$  nonlinear equations and an  $n$ -dimensional minimization problem. To overcome this inconvenience a class of methods has recently attracted much attention (Ref. 5,6,7,8,9): it can be viewed as a family of methods to solve the system of  $(n+m)$  equations arising from the first order optimality conditions for problem (1.1):

$$(1.5) \quad \nabla \ell(x^*, \lambda^*) = \begin{bmatrix} \nabla_x \ell(x^*, \lambda^*) \\ c(x^*) \end{bmatrix} = 0;$$

Notice that the augmented Lagrangian  $\ell(x, \lambda, r)$  could also be used. It is possible to express both Newton and Quasi-Newton iterations for the solution of (1.5) as

$$(1.6) \quad x^{k+1} = x^k + t_k d^k$$

where the direction  $d^k$  is solution of the quadratic programming problem (Q.P.)

$$(1.7) \quad \text{Min } \left\{ \frac{1}{2} \langle d, M_k d \rangle + \langle \nabla f(x^k), d \rangle \mid d \in \mathbb{R}^n \text{ s.t. } c(x^k) + [\nabla c(x^k)]^T d = 0 \right\};$$

$M_k$  is the Hessian with respect to  $x$  of the Lagrangian function or an approximation based upon a Quasi-Newton update formula. The stepsize  $t_k$  is introduced to guarantee the global convergence of the method; it can be selected to achieve at each iteration a sufficient decrease of the (non-differentiable) exact penalty function

$$(1.8) \quad \Phi(x, r) = f(x) + r \sum_{i=1}^m |c_i(x)|,$$

as suggested by Han (Ref.10). If, after a finite number of iterations, the stepsize  $t_k$  can be chosen equal to 1, then the Quasi-Newton method (1.6) has a superlinear rate of convergence (Refs. 8,9). Notice that in contrast with the superlinearly convergent Quasi-Newton method presented in (Ref.1), Han-Powell method requires the updating of an  $n$ -dimensional positive definite matrix  $M_k$ . Although such Quasi-Newton methods involve a sequence of solutions of Q.P., their advantage over the primal methods and multipliers methods lies in the key property that the Q.P. can be solved efficiently in a finite number of iterations.

In this paper we propose a Quasi-Newton method which offers the nice features of both the Han-Powell method and the Quasi-Newton method of (Ref.1). It consists of a sequence of quadratic programs (1.7) where the matrix  $M_k$  is of rank  $(n-m)$  and can be expressed and updated in terms of an  $(n-m)$ -dimensional positive definite symmetric matrix  $H_k$  and its update  $H_{k+1}$ , defined by the generalized Broyden-Fletcher-Collfabb-Shanno formula of (Ref.1)

$$(1.9) \quad H_{k+1} = \left( I - \frac{s^k (y^k)^T}{\langle y^k, s^k \rangle} \right) H_k \left( I - \frac{y^k (s^k)^T}{\langle y^k, s^k \rangle} \right) + \frac{s^k (s^k)^T}{\langle y^k, s^k \rangle},$$

where

$$\begin{aligned} y^k &= g^{k+1} - g^k, \\ s^k &= -t_k H_k g^k, \end{aligned}$$

$g^k$  and  $g^{k+1}$  being the reduced gradients defined now at the non-feasible points  $x^k$  and  $x^{k+1}$ . This new method can thus be interpreted as a generalization of the Quasi-Newton method of (Ref.1) to non-feasible points and with partial restoration (one step of the modified Newton's method for the solution of the constraint equations). Geometrically each iteration can be viewed as a combination of a step in the tangent space to the manifold  $C^k = c^{-1}(c(x^k))$ , "parallel" to  $C = c^{-1}(0)$ , and a partial restoration step to improve (and no more enforce) feasibility. This method presents some similarities with the combined gradient-restoration algorithm experienced by Miele et al. (Ref.11) and with the Quasi-Newton Feasibility Improving GRG method mentioned by Tanabe (Ref.12).

In section 2 we state regularity assumptions on problem (1.1) and review the optimality conditions in terms of the ordinary and augmented Lagrangian functionals. A family of Quasi-Newton directions is presented in Section 3 as solutions of quadratic programming problems. We also show the relation between such Quasi-Newton methods and multipliers methods. The rate of convergence of algorithm (1.6) (with  $t_k = 1$ ) is investigated in Section 4 ; the two-steps superlinear convergence result of Powell (Ref.9) is shown to hold for the Quasi-Newton method defined by a sequence of matrices  $M_k$  of rank  $(n-m)$ . Superlinear convergence is also established for a simply modified algorithm. Finally in Section 5 we review methods for selecting the stepsize.

## 2. ASSUMPTIONS AND OPTIMALITY CONDITIONS

We consider the nonlinear programming problem

$$(2.1) \quad \text{Min } \{f(x) \mid x \in \mathbb{R}^n \text{ } c(x) = 0\}$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  ( $m \leq n$ ) are  $\mathcal{C}^\sigma$  differentiable functions ( $\sigma \geq 2$ ). In addition we assume that the map  $c$  is a submersion, i.e. for all  $x \in \mathbb{R}^n$  the Jacobian map  $c'(x) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$  is of full rank  $m$ . The regular value theorem (see Milnor (Ref.13)) implies that the set  $C_w = c^{-1}(w)$  is a differential manifold of class  $\mathcal{C}^\sigma$  for all  $w \in \mathbb{R}^m$ ; thus any point  $x \in \mathbb{R}^n$  belongs to a differential manifold  $C_{c(x)}$  and it is possible to define  $T_x$  the tangent space to  $C_{c(x)}$  at  $x$  :

$$(2.2) \quad T_x = \{y \in \mathbb{R}^n \mid c'(x)(y) = 0\}.$$

Remark 2.1. This assumption on  $c$  is stronger than the regularity assumption introduced in (Ref.1), namely 0 is a regular value of  $c$ . Under the last assumption there may exist points  $x \in \mathbb{R}^n - c^{-1}(0)$  where the Jacobian map  $c'(x)$  is not surjective; such points are called critical points of  $c$ . However, if  $c$  is a  $\mathcal{C}^\sigma$  differentiable map with

$$\sigma > n - m,$$

the Morse-Sard theorem (see e.g. (Ref.14)) establishes that the image by  $c$  of the critical points is a set of measure zero in  $\mathbb{R}^m$ ; hence  $C_w$  is a differential manifold for "almost" every  $w \in \mathbb{R}^m$ . ■

We recall that it is possible to characterize a solution  $x^*$  of problem (2.1) in terms of the Lagrangian functional  $\ell : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$  defined by

$$(2.3) \quad \ell(x, \lambda) = f(x) + \langle \lambda, c(x) \rangle,$$

(see e.g. Luenberger (Ref.15) for the proofs of the following classical results).

PROPOSITION 2.1. (First-order optimality condition). If  $x^*$  is a local minimum of (2.1) there exists a vector of Lagrange multipliers  $\lambda^* \in \mathbb{R}^m$  such that  $(x^*, \lambda^*)$  is a critical point of the Lagrangian functional (2.3). ■

An equivalent formulation of this condition characterizes  $(x^*, \lambda^*)$  as a solution of the system of  $(n+m)$  equations :

$$(2.4a) \quad \nabla_x \ell(x, \lambda) = \nabla f(x) + \nabla c(x) \cdot \lambda = 0$$

$$(2.4b) \quad \nabla_\lambda \ell(x, \lambda) = c(x) = 0$$

Assuming that the functions  $f$  and  $c$  are  $C^\sigma$  differentiable with  $\sigma \geq 2$ , it is possible to define the Hessian with respect to  $x$  of the Lagrangian functional

$$(2.5) \quad \nabla_{xx}^2 \ell(x, \lambda) = \nabla^2 f(x) + \sum_{i=1}^m \lambda_i \nabla^2 c_i(x)$$

and derive second-order characterizations of a solution  $x^*$ .

PROPOSITION 2.2. (Second-order necessary optimality condition). If  $x^*$  is a local minimum of (2.1) there exists a  $\lambda^* \in \mathbb{R}^m$  such that  $(x^*, \lambda^*)$  is a critical point of  $\ell(x, \lambda)$  and the restriction of the Hessian form  $\nabla_{xx}^2 \ell(x^*, \lambda^*)$  to the tangent space  $T_{x^*}$  is positive semi-definite. ■

PROPOSITION 2.3. (Second-order sufficient optimality condition). Assume that  $(x^*, \lambda^*)$  is a critical point of the Lagrangian functional and that

$$(2.6) \quad \langle v, L(x^*, \lambda^*) v \rangle > 0 \quad \text{for all} \quad v \in T_{x^*}, v \neq 0,$$

(where  $L(x^*, \lambda^*)$  denotes the  $n \times n$  symmetric matrix of the quadratic form  $\nabla_{xx}^2 \ell(x^*, \lambda^*)$ ). Then  $x^*$  is an isolated local minimum of  $f$  on the manifold  $C = c^{-1}(0)$ . ■

Condition (2.6) allows to distinguish among the critical pairs  $(x^*, \lambda^*)$  which candidates  $x^*$  are actually isolated local minimum of problem



(2.1). Notice that the difference between the necessary and sufficient second-order optimality conditions vanishes if the Lagrangian functional has no degenerated critical points ; this is the case if the restriction of  $f$  to the submanifold  $C$  is a Morse function (see (Ref.1)).

We now establish a result which will be used several times in our analysis.

PROPOSITION 2.4. Assume that  $c$  is a submersion and that the restriction to the tangent space  $T_x$  of the  $n \times n$  matrix  $M_x$  is positive definite. Then, the  $(n+m) \times (n+m)$  matrix  $D(x, M_x)$  defined by

$$(2.7) \quad D(x, M_x) = \begin{bmatrix} M_x & A_x^T \\ A_x & 0 \end{bmatrix},$$

where  $A_x$  denotes the Jacobian matrix of the map  $c$  at  $x$ , is non-singular. ■

Proof : Let  $z = (y, \mu) \in \mathbb{R}^{n+m}$  such that  $D(x, M_x).z = 0$ . By (2.7) we have

$$(2.8) \quad M_x y + A_x^T \mu = 0,$$

$$(2.9) \quad A_x y = 0.$$

Equation (2.9) implies that  $y \in T_x$  and (2.8) yields

$$(2.10) \quad \begin{aligned} \langle y, M_x y \rangle &= -\langle y, A_x^T \mu \rangle \\ &= -\langle A_x y, \mu \rangle = 0. \end{aligned}$$

Since  $M_x$  is positive definite on  $T_x$ , (2.10) implies  $y = 0$  and (2.8) reduces to

$$(2.11) \quad A_x^T \mu = 0,$$

which has for unique solution  $\mu = 0$  since  $A_x$  is of full rank  $m$ . ■

COROLLARY 2.5. Assume that  $(x^*, \lambda^*)$  satisfies the second-order sufficient optimality condition. Then the Hessian matrix of  $\ell(x, \lambda)$  with respect to  $x$  and  $\lambda$  at  $(x^*, \lambda^*)$  is non-singular. ■

Proof : The Hessian matrix of  $\ell(x, \lambda)$  with respect to  $x$  and  $\lambda$  is

$$(2.12) \quad \begin{bmatrix} L(x, \lambda) & A_x^T \\ A_x & 0 \end{bmatrix} = D(x, L(x, \lambda)) \quad ;$$

the result follows from Proposition 2.4 since  $c$  is a submersion by assumption. ■

We conclude this section by a review of optimality conditions in terms of the augmented Lagrangian  $\ell : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^t \rightarrow \mathbb{R}$  defined by

$$(2.13) \quad \ell(x, \lambda, r) = \ell(x, \lambda) + \frac{r}{2} \|c(x)\|^2.$$

If  $(x^*, \lambda^*)$  is a critical point of the ordinary Lagrangian  $\ell(x, \lambda)$  it is also a critical point of the augmented Lagrangian (with respect to  $x$  and  $\lambda$ ) for any  $r > 0$  since

$$(2.14a) \quad \nabla_x \ell(x^*, \lambda^*, r) = \nabla_x \ell(x^*, \lambda^* + rc(x^*)) = \nabla_x \ell(x^*, \lambda^*) = 0 \quad ,$$

$$(2.14b) \quad \nabla_\lambda \ell(x^*, \lambda^*, r) = \nabla_\lambda \ell(x^*, \lambda^*) = 0 \quad .$$

The Hessian of the augmented Lagrangian is given by

$$(2.15) \quad \nabla_{xx}^2 \ell(x, \lambda, r) = \nabla_{xx}^2 \ell(x, \lambda + rc(x)) + r \nabla c(x) \cdot [\nabla c(x)]^T \quad ;$$

hence, denoting by  $L_r(x, \lambda)$  the matrix of the quadratic form  $\nabla_{xx}^2 \ell(x, y, r)$ ,

$$\langle v, L_r(x, \lambda)v \rangle = \langle v, L(x, \lambda + rc(x))v \rangle \quad \text{for all } v \in T_x$$

and the second-order optimality conditions (Propositions 2.2 and 2.3) can be equivalently expressed in term of the augmented Lagrangian functional (2.13).

The key benefit of introducing this augmented functional lies in the following now classical result (see e.g. (Ref.4)).

PROPOSITION 2.6. Assume that  $(x^*, \lambda^*)$  satisfies the second-order sufficient optimality condition. Then there exists a scalar  $r^* \geq 0$  such that the  $(n \times n)$  symmetric matrix  $L_r(x^*, \lambda^*)$  is positive definite for all  $r \geq r^*$ . ■

The continuity of  $\nabla^2 f(x)$  and  $\nabla^2 c_i(x)$  for  $i = 1, \dots, m$  guarantees that the matrix  $L_r(x, \lambda)$  remains positive definite in a neighborhood  $B(x^*, \epsilon_1) \times B(\lambda^*, \epsilon_2)$  of the critical point  $(x^*, \lambda^*)$  which can thus be characterized as a (local) saddle-point of the (locally) convex-concave function  $\ell(x, y, r)$ . This is the basis for the interpretation of the multipliers method as a duality scheme. In particular we have

$$(2.16) \quad \ell(x^*, \lambda^*, r) = \min_{x \in B(x^*, \epsilon_1)} \ell(x, \lambda^*, r) \quad \text{for all } r \geq r^* .$$

A somewhat similar result holds for the (non-differentiable) exact penalty function  $\Phi : \mathbb{R}^n \times \mathbb{R}^+ \rightarrow \mathbb{R}$  defined by

$$(2.17) \quad \Phi(x, r) = f(x) + r \sum_{i=1}^m |c_i(x)|$$

Let  $(x^*, \lambda^*)$  satisfies the second-order sufficient optimality condition and let  $d \in T_{x^*}$ , sufficiently small, such that

$$\ell(x^* + d, \lambda^*) = \ell(x^*, \lambda^*) + \langle d, L(x^* + \theta d, \lambda^*) d \rangle \geq \ell(x^*, \lambda^*) .$$

Noticing that

$$\Phi(x^*, r) = f(x^*) = \ell(x^*, \lambda^*) ,$$

we have

$$\begin{aligned} \Phi(x^* + d, r) &= \ell(x^* + d, \lambda^*) + \sum_{i=1}^m [r |c_i(x^* + d)| - \lambda_i^* c_i(x^* + d)] \\ &\geq \Phi(x^*, r) , \end{aligned}$$

provided

$$(2.18) \quad r \geq r^* = \max_i |\lambda_i^*| ;$$

hence

$$(2.19) \quad \Phi(x^*, r) = \min_{d \in B(0, \varepsilon_1) \cap T_x^*} \Phi(x^* + d, r) \quad \text{for all } r \geq r^* .$$

### 3. A CLASS OF QUASI-NEWTON DIRECTIONS

#### 3.1. A Quadratic Programming Problem.

The family of Newton and Quasi-Newton methods for solving the system of  $(n+m)$  equations arising from the first-order optimality conditions (2.4) generates iteratively a sequence of approximate solutions  $\{(x^k, \lambda^k)\}$  according to

$$(3.1) \quad \begin{bmatrix} M_k & A_k^T \\ A_k & 0 \end{bmatrix} \begin{bmatrix} x^{k+1} - x^k \\ \lambda^{k+1} - \lambda^k \end{bmatrix} + \begin{bmatrix} \nabla_x \ell(x^k, \lambda^k) \\ c(x^k) \end{bmatrix} = 0 ,$$

where  $A_k$  denotes the Jacobian matrix of the map  $c$  at  $x^k$  and the  $n \times n$  matrix  $M_k$  is the Hessian matrix of the Lagrangian with respect to  $x$ ,  $L(x^k, \lambda^k)$ , (in Newton's method) or an approximation of it (in Quasi-Newton methods). The choice  $M_k = L(x^k, \lambda^k)$ , which requires the computation of second derivatives of  $f$  and  $c$ , has first been presented by Wilson (Ref.16) and the convergence properties of the method established by Robinson (Ref.6). To avoid the need for second-order derivatives Garcia-Palomares and Mangasarian (Ref.7) have introduced an approximation  $M_k$  of  $L(x^k, \lambda^k)$  obtained by updating a large  $(n+m) \times (n+m)$  matrix at each iteration and requiring only first-order information. The wasteful character of this procedure has been evidenced by Han (Ref.8) who proposed to update directly the  $n \times n$  matrix  $M_k$  by a Quasi-Newton update formula.

Our aim in this paper is to go a step further, namely to show that it is only necessary to update a reduced matrix of dimension  $(n-m) \times (n-m)$ , as in the Quasi-Newton methods along geodesics proposed in (Ref.1). Before we present our method, we provide a general analysis of method (3.1) and show the convenience of requiring that the restriction of  $M_k$  to the tangent space  $T_k$  at  $x^k$  be positive definite.

Using the notation

$$(3.2) \quad d^k = x^{k+1} - x^k ,$$

the definition (3.1) of the Quasi-Newton method yields the system of equations

$$(3.3a) \quad \nabla f(x^k) + M_k d^k + A_k^T \lambda^{k+1} = 0 ,$$

$$(3.3b) \quad c(x^k) + A_k d^k = 0 ,$$

which characterizes  $(d^k, \lambda^{k+1})$  as a solution\* of the quadratic programming problem (Q.P.)

$$(3.4) \quad \text{Min } \{ \langle \nabla f(x^k), d \rangle + \frac{1}{2} \langle d, M_k d \rangle \mid d \in \mathbb{R}^n \text{ s.t. } c(x^k) + A_k d = 0 \} ;$$

Thus method (3.1) can be viewed as consisting of solving a sequence of Q.P. defined recursively ; since such problems can be efficiently solved in a finite number of iterations, this approach presents a clear advantage over primal methods, which require successive iterative solutions of systems of nonlinear equations, as well as over multiplier methods, which require successive minimizations of the augmented Lagrangian. Powell (Ref. 9) has explored further this method and shown that difficulties may arise when the Q.P. has several solutions. However, in the present case, we can simply overcome this difficulty.

PROPOSITION 3.1. If the restriction of the matrix  $M_k$  to the tangent space  $T_k$  at  $x^k$  is positive definite, the Q.P. (3.4) has a unique solution. ■

Proof : The first-order optimality conditions (3.3) for problem (3.4) form a system of linear equations in  $(d, \lambda)$  of matrix  $D(x^k, M_k)$  defined by (2.7). Proposition 2.4 establishes that, under our assumptions,  $D(x^k, M_k)$  is non singular ; hence (3.3) admits a unique solution  $(d^k, \lambda^{k+1})$ . The Hessian of the Lagrangian functional associated to (3.4) is simply the matrix  $M_k$  ; hence  $(d^k, \lambda^{k+1})$  satisfy also the second-order sufficient optimality condition and  $d^k$  actually achieves the minimum of the quadratic programming problem (3.4). ■

---

\* By a solution of the Q.P. we mean the couple consisting of a minimizing element and of the corresponding vector of Lagrange multipliers associated to the constraints.

### 3.2. A change of coordinates

In (Ref.1) we have found convenient to associate to the full rank Jacobian map  $c'(x) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$  of  $c$  at  $x$  of matrix  $A_x$  a linear map of  $\mathcal{L}(\mathbb{R}^n, \mathbb{R}^{n-m})$  defined by an  $(n-m) \times n$  matrix  $Z_x$  such that

$$\eta(Z_x) \cap \eta(A_x) = 0 ;$$

the  $n \times n$  matrix  $S_x$  defined by

$$(3.5) \quad S_x = \begin{bmatrix} A_x \\ Z_x \end{bmatrix}$$

is non-singular and thus defines a change of coordinate in  $\mathbb{R}^n$ . We recall the following result established in (Ref.1, Prop. 2.1.).

PROPOSITION 3.2. Assume that  $c$  is a submersion. Let  $A_x^-$  be a right inverse for  $A_x$  (i.e.  $A_x \cdot A_x^- = I_m$ ). Then there exists a unique matrix  $Z_x$  of full rank  $(n-m)$  with right inverse  $Z_x^-$  satisfying

$$(3.6) \quad Z_x \cdot A_x^- = 0 \quad , \quad A_x \cdot Z_x^- = 0.$$

The matrix  $S_x$  defined by (2.5) is non-singular and its inverse is given by

$$(3.7) \quad S_x^{-1} = [A_x^-, Z_x^-] .$$

If  $c$  is a  $\mathcal{C}^\sigma$  map,  $A_x^-$ ,  $Z_x$ ,  $Z_x^-$  can be chosen as  $\mathcal{C}^{\sigma-1}$  differentiable functions of  $x$  in  $\mathbb{R}^n$ . ■

This change of coordinate allows us to express the unique solution of the Q.P. (3.4) in a form convenient for computations.

PROPOSITION 3.3. Assume that the restriction of the matrix  $M$  to the  $(n-m)$  dimension subspace  $\{y \in \mathbb{R}^n \mid Ay = 0\}$  is positive definite. Then

$$(3.8) \quad \begin{bmatrix} M & A^T \\ A & 0 \end{bmatrix}^{-1} = \begin{bmatrix} Z^{-T} H Z^{-T} & (I - Z^{-T} H Z^{-T} M) A^{-T} \\ A^{-T} (I - M Z^{-T} H Z^{-T}) & -A^{-T} (I - M Z^{-T} H Z^{-T}) M A^{-T} \end{bmatrix}$$

where  $H = (Z^{-T} M Z^{-T})^{-1}$ . ■

Proof : We want to solve in  $(y, \mu)$  the system of equations

$$(3.9a) \quad M y + A^T \mu = g,$$

$$(3.9b) \quad A y = h;$$

it is non-singular by Proposition 3.1 and thus has a unique solution. Using the change of coordinate defined by  $S_x^{-1}$  we can write

$$(3.10) \quad y = A^{-T} w + Z^{-T} z$$

where  $w = A y \in \mathbb{R}^m$  and  $z = Z y \in \mathbb{R}^{n-m}$ . Equations (3.6) (3.9b) and (3.10) imply that  $w = h$ ; premultiplying (3.9a) by the non-singular matrix  $(S_x^{-1})^T$  yields the system of  $n$  equations in  $z, \mu$ :

$$(3.11a) \quad \mu + A^{-T} M Z^{-T} z + A^{-T} M A^{-T} h = A^{-T} g$$

$$(3.11b) \quad Z^{-T} M Z^{-T} z + Z^{-T} M A^{-T} h = Z^{-T} g.$$

The  $(n-m) \times (n-m)$  matrix  $Z^{-T} M Z^{-T}$  represents the restriction to the subspace  $\{y \in \mathbb{R}^n \mid A y = 0\} = \{Z^{-T} z \mid z \in \mathbb{R}^{n-m}\}$ ; by assumption it is positive definite, hence invertible. Let  $H = (Z^{-T} M Z^{-T})^{-1}$ ;  $H$  is positive definite. Solving (3.11b) for  $z$  and getting  $\mu$  from (3.11a) after substitution, we obtain formula (3.8). ■

This formula generalizes the well-known formula for the inverse in partitioned form (see e.g. (Ref.17)) valid only when  $M$  is non-singular. It includes as a special case a formula used by Powell (Ref.18), Tapia



(Ref.19) valid only when  $M$  is non-singular. A more complex formula is given by Tanabe (Ref.12) in term of a generalized inverse  $A^-$  of  $A$  (such that  $A A^- A = A$ ) but is of little value in practical computation ; assuming that  $A$  is of full rank the generalized inverse is also a right inverse and Tanabe's formula should reduce to (3.8).

### 3.3. The Quasi-Newton Direction

It is legitimate at this point to define, as in (Ref.1), the reduced gradient of  $f$  at the nonfeasible point  $x$  as the  $(n-m)$  dimensional vector

$$(3.12) \quad g(x) = Z_x^{-T} \nabla f(x) ;$$

obviously the reduced gradient depends upon the change of coordinate defined by  $S_x$ . The Quasi-Newton direction defined as the solution  $d^k$  of the Q.P. (3.4) can be expressed using (3.8) as

$$(3.13) \quad d^k = - Z_k^- H_k (g^k - Z_k^{-T} M_k A_k^- c^k) - A_k^- c^k ,$$

where the subscripts and superscripts  $k$  indicate that the respective matrices and vectors are evaluated at the current iterate  $x^k$ . Formula (3.13) indicates that the Quasi-Newton direction  $d^k$  is a combination of a direction in the tangent space  $T_k$  to the submanifold  $C_k = c^{-1}(c^k)$  at  $x^k$  and a direction pointing toward the constraint manifold  $C = c^{-1}(0)$  (since  $-A_k^- c^k$  can be viewed as the first step of a Newton's method starting from  $x^k$  to solve the system of equations  $c(x) = 0$ ).

Formula (3.13) defines a double family of Quasi-Newton directions. On one hand it depends upon the choice of the change of coordinates  $S_x$ , i.e. of the choice of the right inverse  $A_x^-$  of  $A_x$  which then conditions the choices of  $Z_x$  and  $Z_x^-$  according to Proposition 3.2. Two convenient choices have been considered in (Ref.1). The partitioned right inverse consists in partitioning of  $A_x$  as

$$(3.14) \quad A_x = [B, D] ,$$

(where B is an  $m \times m$  non singular matrix) and choosing respectively

$$(3.15) \quad A_x^- = \begin{bmatrix} B^{-1} \\ 0 \end{bmatrix} , \quad Z_x = [0, I_{n-m}] , \quad Z_x^- = \begin{bmatrix} -B^{-1} D \\ I_m \end{bmatrix} .$$

The QR right inverse consists in factorizing  $A_x$  as

$$(3.16) \quad A_x = [L, 0] \begin{bmatrix} Q_1 \\ Q_2 \end{bmatrix} ,$$

(where L is an  $m \times m$  lower triangular matrix and  $Q_1$  and  $Q_2$  are submatrices of respective dimensions  $m$  and  $(n-m)$  of an orthogonal matrix) and choosing respectively

$$(3.17) \quad A_x^- = Q_1^T L^{-1} \quad Z_x = Q_2 \quad Z_x^- = Q_2^T .$$

On the other hand the Quasi Newton direction (3.13) depends upon the scheme adopted to update the matrix  $M_k$ . Many update formulae are available since we do not require  $M_k$  to be symmetric (in fact the quadratic programming problem (3.4) can be defined for a non symmetric  $M_k$  by substituting  $M_k$  by  $\frac{1}{2}(M_k + M_k^T)$  in its formulation, as noticed by Han (Ref.10)). However our analysis requires that  $M_k$  remains positive definite on the (changing) subspace  $T_k$ . A convenient scheme to achieve this property in a very simple way consists in considering matrices  $M_k$  satisfying

$$(3.18) \quad M_k A_k^- = 0 ;$$

in this case the Quasi-Newton direction is given by

$$(3.19) \quad d^k = - Z_k^- H_k g^k - A_k^- c^k ,$$

while the corresponding Lagrange multipliers vector of the Q.P. (3.4) is

$$(3.20) \quad \lambda^{k+1} = - A_k^{-T} \nabla f(x^k) .$$

Observe that the solution of (3.4) depends only upon  $H_k = (Z_k^{-T} M_k Z_k^{-1})^{-1}$  and it is thus only necessary to update  $G_k = H_k^{-1}$  as a positive definite symmetric matrix of dimension  $(n-m)$  approximating the restriction to the tangent space of the Hessian of the Lagrangian.

By analogy with the reduced Quasi-Newton method of (Ref.1) we can think of updating  $G_k$  by the Generalized Broyden-Fletcher-Goldfarb-Shanno formula

$$(3.21) \quad G_{k+1} = G_k + \frac{y_k (y_k)^T}{\langle y^k, s^k \rangle} - \frac{G_k s^k (s^k)^T G_k}{\langle s^k, G_k s^k \rangle}$$

where

$$(3.22a) \quad s^k = x^{k+1} - x^k ,$$

$$(3.22b) \quad y^k = g^{k+1} - g^k ;$$

(we introduce the notation  $s^k$  to allow for the possibility of introducing a stepsize  $t_k$  in the iteration  $x^{k+1} = x^k + t_k d^k$ ).

However there is no guarantee in the new method that

$$\langle y^k, s^k \rangle > 0$$

which is a necessary and sufficient condition for  $G_{k+1}$  to be positive definite if  $G_k$  is. We must therefore modify the update formula (3.21) using a device suggested by Powell (Ref.20) :

$$(3.23) \quad G_{k+1} = G_k + \frac{z^k (z^k)^T}{\langle z^k, s^k \rangle} - \frac{G_k s^k (s^k)^T G_k}{\langle s^k, G_k s^k \rangle} ,$$

where

$$(3.24a) \quad z^k = \theta_k y^k + (1-\theta_k) G_k s^k$$

and  $\theta_k$  is a scalar between 0 and 1 chosen according to

$$(3.24b) \quad \theta_k = \begin{cases} 1 & \text{if } \langle y^k, s^k \rangle \geq \sigma \langle s^k, G_k s^k \rangle, \\ \frac{(1-\sigma) \langle s^k, G_k s^k \rangle}{\langle s^k, G_k s^k \rangle - \langle y^k, s^k \rangle} & \text{otherwise,} \end{cases}$$

with  $\sigma \in (0, \frac{1}{2})$  (Powell suggests to use  $\sigma = 0.2$ ).

Notice that  $H_k = G_k^{-1}$  can then be updated directly according to

$$(3.25) \quad H_{k+1} = \left( I - \frac{s^k (z^k)^T}{\langle z^k, s^k \rangle} \right) H_k \left( I - \frac{z^k (s^k)^T}{\langle z^k, s^k \rangle} \right) + \frac{s^k (s^k)^T}{\langle z^k, s^k \rangle},$$

although we do not recommend the use of this formula in practical computations for its potential numerical instability.

#### 3.4. Relation with multipliers methods:

We could also define Newton and Quasi-Newton methods for solving the system of  $(n+m)$  equations arising from the first-order optimality conditions expressed in term of the augmented Lagrangian.

$$(3.26a) \quad \nabla_x \ell(x, \lambda, r) = \nabla_x \ell(x, \lambda) + r \nabla c(x) \quad c(x) = 0,$$

$$(3.26b) \quad \nabla_\lambda \ell(x, \lambda, r) = c(x) = 0.$$

This scheme leads to the iterative definition of a sequence  $\{(y^k, \mu^k)\}$  according to

$$(3.27) \quad \begin{bmatrix} N_k & A_k^T \\ A_k & 0 \end{bmatrix} \begin{bmatrix} y^{k+1} - y^k \\ \mu^{k+1} - \mu^k \end{bmatrix} + \begin{bmatrix} \nabla_x \ell(y^k, \mu^k) + r A_k^T c(y^k) \\ c(y^k) \end{bmatrix} = 0,$$

where  $N_k$  is an approximation of the Hessian matrix of the augmented Lagrangian given by (2.15); it is thus legitimate to define

$$(3.28) \quad N_k = M_k + r A_k^T A_k$$

where  $M_k$  is an approximation of the Hessian matrix of the ordinary Lagrangian. Formula (3.8) can again be used to compute the solution of (3.27). Noticing that (3.6) implies

$$Z_k^{-T} N_k Z_k^- = Z_k^{-T} M_k Z_k^- = H_k^{-1} ,$$

it is easy to show that, starting from  $(y^k, \mu^k) \equiv (x^k, \lambda^k)$ , method (3.27) generates the same iterate  $x^{k+1}$  as method (3.1), i.e.

$$(3.29) \quad y^{k+1} - y^k = d^k = - Z_k^- H_k (g^k - Z_k^{-T} M_k A_k^- c^k) - A_k^- c^k ,$$

while the Lagrange multipliers  $\mu^{k+1}$  can be deduced from  $\lambda^{k+1}$ , solution of (3.4), by

$$(3.30) \quad \mu^{k+1} = \lambda^{k+1} + r c^k ,$$

a result established in the particular case where  $M_k$  is non-singular by Powell (Ref.18) and Tapia (Ref.19).

Suppose that  $(y^k, \lambda^k)$  is in a small enough neighborhood of  $(x^*, \lambda^*)$  satisfying the second-order sufficient optimality condition and that  $r$  is chosen large enough so that  $N_k$  given by (3.28) is positive definite, hence invertible (this is possible by Proposition 2.6). The first block of equations of (3.27) can be solved for  $(y^{k+1} - y^k)$  and yields

$$(3.31) \quad y^{k+1} = y^k - N_k^{-1} (\nabla f(y^k) + \nabla c(y^k) (\mu^{k+1} + r c(y^k))) ;$$

equation (3.31) can be viewed as one step of a Quasi-Newton method starting from  $y^k$  to solve the unconstrained (locally convex) minimization problem

$$(3.32) \quad \min_{y \in B(x^*, \varepsilon_1)} \ell(y, \mu^{k+1}, r) .$$

Notice that the multipliers vector  $\mu^{k+1}$  is then given by

$$(3.33) \quad \mu^{k+1} = \mu^k + (A_k N_k^{-1} A_k^T)^{-1} (c^k - A_k N_k^{-1} \nabla_x \ell(y^k, \mu^k, r)) ;$$

If the minimization phase (3.32) had been performed exactly at the previous iteration, we would have  $\nabla_x \ell(y^k, \mu^k, r) = 0$  and (3.33) would reduce to

$$(3.34) \quad \mu^{k+1} = \mu^k + (A_k N_k^{-1} A_k^T)^{-1} c^k.$$

Iteration (3.34) can be viewed as the  $k$ -th step of a Quasi-Newton to maximize the dual functional  $\Psi_r : B(\lambda^*, \epsilon_2) \rightarrow \mathbb{R}$  associated to the augmented Lagrangian

$$(3.35) \quad \Psi_r(\mu) \triangleq \min_{y \in B(y^*, \epsilon_1)} \ell(y, \mu, r),$$

since  $-A_k N_k^{-1} A_k^T$  is an approximation to  $\nabla^2 \Psi_r(\mu^k) = -A_k (L_r(y^k, \mu^k))^{-1} A_k^T$  (see e.g. (Ref.4)).

Thus method (3.27), which is related to our general Quasi-Newton method (3.4) as exhibited by (3.29) (3.30), can be interpreted as a particular efficient implementation of a Quasi-Newton method for solving the dual problem  $\max \Psi_r(\mu)$ , where the minimization phase (3.32) is performed only approximately by one step of a Quasi-Newton method. Such a method has been called by Tapia (Ref.19) a diagonalized multipliers method; a particular implementation (corresponding to  $M_k = I$ ) has been experimented by Miele et al. (Ref.21). See also Tapia (Ref.22) for a related discussion.

#### 4. SUPERLINEAR CONVERGENCE

The object of this section is to show that the greatly simplifying strategy (3.19), (3.20) requiring only the updating of a reduced  $(n-m) \times (n-m)$  matrix  $G_k$  (approximating the restriction to the tangent space of the Hessian of the Lagrangian) still preserves the attractive superlinear rate of convergence of the Quasi-Newton method of Han-Powell requiring the update of the  $n \times n$  matrix  $M_k$ .

##### 4.1. Two-steps Superlinear Convergence

We consider the iterative method

$$(4.1) \quad x^{k+1} = x^k + d^k,$$

where  $d^k$  is the general Quasi-Newton direction defined in Section 3.3

$$(4.2) \quad d^k = -Z_k^{-1} H_k (g^k - Z_k^{-T} M_k A_k^{-1} c^k) - A_k^{-1} c^k$$

with

$$(4.3a) \quad H_k = (Z_k^{-T} M_k Z_k^{-1})^{-1},$$

$$(4.3b) \quad g^k = Z_k^{-T} \nabla f(x^k),$$

$$(4.3c) \quad c^k = c(x^k).$$

Assume that the sequence  $\{x^k\}$  converges to a local minimum  $x^*$  of  $f$  on  $C$  and let  $\lambda^*$  be the corresponding Lagrange multipliers vector associated to the constraint equations. Assume moreover that  $(x^*, \lambda^*)$  satisfies the second-order optimality condition.

Assume also that the method uses a sequence of bounded matrices  $M_k$  such that

$$(4.4) \quad |\langle v, M_k v \rangle| \leq m_1 \|v\|^2 \quad \text{for all } v \in \mathbb{R}^n \text{ and all } k$$

and that the matrices  $H_k$  remain positive definite and such that

$$(4.5) \quad m_2 \|p\|^2 \leq \langle p, H_k p \rangle \leq m_3 \|p\|^2 \quad \text{for all } p \in \mathbb{R}^{n-m} \text{ and all } k.$$

The direction  $d^k$  is solution of the Q.P. (3.4) to which is associated a Lagrange multipliers vector  $\lambda^{k+1} \in \mathbb{R}^m$  given (by application of (3.8)) by

$$(4.6) \quad \lambda^{k+1} = -A_k^{-T} (I - M_k Z_k^{-1} H_k Z_k^{-T}) (\nabla f(x^k) - M_k A_k^{-1} c(x^k)).$$

PROPOSITION 4.1. If  $\{x^k\} \rightarrow x^*$  the sequence of Lagrange multipliers  $\{\lambda^k\}$  constructed by the successive Q.P. converges to  $\lambda^*$ . ■

Proof : Since  $(d^k, \lambda^{k+1})$  is the unique solution of Q.P. (3.4) it satisfies

$$(4.7) \quad M_k d^k + A_k^T \lambda^{k+1} + \nabla f(x^k) = 0.$$

Since  $\{x^k\} \rightarrow x^*$ ,  $d^k \rightarrow 0$ ; the pair  $(x^*, \lambda^*)$  satisfies the first-order optimality condition

$$(4.8) \quad \nabla f(x^*) + \nabla c(x^*) \lambda^* = 0.$$

Subtracting (4.8) from (4.7), we obtain

$$H_k d^k + A_k^T (\lambda^{k+1} - \lambda^*) + \nabla f(x^k) - \nabla f(x^*) + (\nabla c(x^k) - \nabla c(x^*)) \lambda^* = 0.$$

which by continuity of  $\nabla f$  and  $\nabla c$  and (4.4) yields  $\lambda^{k+1} \rightarrow \lambda^*$ . ■

PROPOSITION 4.2. If  $(x^*, \lambda^*)$  satisfies the second-order sufficient optimality condition, there exists  $K_1 > K_2 > 0$  such that

$$(4.9) \quad K_2 [\|g^k\| + \|c^k\|] \leq \|x^k - x^*\| \leq K_1 [\|g^k\| + \|c^k\|]. \quad \blacksquare$$

Proof : By Corollary 2.5 the Hessian matrix  $D(x^*, L(x^*, \lambda^*))$  of the Lagrangian with respect to  $x$  and  $\lambda$  at  $(x^*, \lambda^*)$  is non-singular; by continuity of the second derivatives of  $f$  and  $c$  the matrix  $D(x, L(x, \lambda))$  remains non-singular in a neighborhood  $B(x^*, \epsilon_1) \times B(\lambda^*, \epsilon_2)$ . Suppose that  $x^k \in B(x^*, \epsilon_1)$  and  $\lambda^{k+1} \in B(\lambda^*, \epsilon_2)$ ;



Taylor expansion of  $\nabla \ell(x^k, \lambda^{k+1})$  around the critical point  $(x^*, \lambda^*)$  yields

$$(4.10) \quad \begin{aligned} \nabla \ell(x^k, \lambda^{k+1}) &= \int_0^1 D(x^* + t(x^k - x^*), L(x^* + t(x^k - x^*), \lambda^* + t(\lambda^{k+1} - \lambda^*))) \begin{bmatrix} x^k - x^* \\ \lambda^{k+1} - \lambda^* \end{bmatrix} dt \\ &= \bar{D} \begin{bmatrix} x^k - x^* \\ \lambda^{k+1} - \lambda^* \end{bmatrix}, \end{aligned}$$

where the matrix  $\bar{D}$  is non singular. Hence

$$(4.11) \quad x^k - x^* = (\bar{D}^{-1})_{11} \nabla_x \ell(x^k, \lambda^{k+1}) + (\bar{D}^{-1})_{12} c(x^k)$$

Using the expression (4.6) of  $\lambda^{k+1}$ , we have

$$(4.12) \quad \nabla_x \ell(x^k, \lambda^{k+1}) = M_k Z_k^- H_k Z_k^{-T} g^k + (I - M_k Z_k^- H_k Z_k^{-T}) M_k A_h^- c^k.$$

By proposition 3.2,  $A_x^-$  and  $Z_x^-$  are  $C^{\sigma-1}$  functions of  $x$  hence bounded on  $B(x^*, \epsilon_1)$ ; formulae (4.11) and (4.12) together with (4.4), (4.5) yield (4.9). ■

Under the same assumptions we obtain the estimate of the norm of the Quasi-Newton direction.

PROPOSITION 4.3. There exists  $K_3 > K_4 > 0$  such that

$$(4.14) \quad K_4 [||g^k|| + ||c^k||] \leq ||d^k|| \leq K_3 [||g^k|| + ||c^k||]. \quad \blacksquare$$

We establish a bound on the constraints violation at each iteration.

PROPOSITION 4.4. There exists  $K_5 > 0$  such that

$$(4.15) \quad ||c(x^k + d^k)|| \leq K_5 ||d^k||^2. \quad \blacksquare$$

Proof : Since the functionals  $c_i$  are  $C^\sigma$  differentiable with  $\sigma \geq 2$  we have

$$\gamma_i = \sup_{x \in B(x^*, \epsilon_1)} \left( \sup_{||v||=1} ||c_i''(x) \cdot v|| \right) < +\infty.$$

where  $m_4$  is a bound on  $B(x^*, \epsilon_1)$  of the derivatives of  $L(., \lambda^{k+1})$ . Observe from (4.2) that we can write

$$d^k = -Z_k^{-1} p^k + q^k,$$

where  $p^k = H_k g^k$  and  $q^k = (Z_k^{-1} H_k Z_k^{-T} M_k - I) A_k^{-1} c^k$  is such that

$$(4.19) \quad ||q^k|| \leq m_5 ||c^k||.$$

Combining (4.16), (4.17), (4.18), (4.19) we obtain

$$\begin{aligned} ||g^{k+1}|| &\leq ||Z_k^{-T} \nabla_x \ell(x^{k+1}, \lambda^{k+1})|| + ||(Z_{k+1}^{-1} - Z_k^{-1})^T \nabla_x \ell(x^{k+1}, \lambda^{k+1})|| \\ &\leq ||Z_k^{-T} [L(x^k, \lambda^{k+1}) - M_k] Z_k^{-1} p^k|| + K_6 ||d^k||^2 + K_7 ||c^k|| \end{aligned}$$

PROPOSITION 4.5. Assume that  $f$  and  $c$  are  $C^\sigma$  differentiable with  $\sigma \geq 3^*$ ; then there exists  $K_6, K_7 > 0$  such that

$$(4.20) \quad ||g^{k+1}|| \leq ||[Z_k^{-T} L(x^k, \lambda^{k+1}) Z_k^{-1} - H_k^{-1}] p^k|| + K_6 ||d^k||^2 + K_7 ||c^k||;$$

with  $p^k = -H_k g^k$ . ■

We are now in position to establish the two steps superlinear convergence for the class of methods (4.1) (4.2), extending a result established by Powell (Ref.9) for a particular implementation.

THEOREM 4.1. Assume that  $f$  and  $c$  are  $C^\sigma$  differentiable with  $\sigma \geq 3^*$  and that the general Quasi-Newton method (4.1) (4.2) generates a sequence of approximate solutions  $\{x^k\}$  together with Lagrange multipliers  $\{\lambda^k\}$  (given by (4.6)) such that  $\{(x^k, \lambda^k)\}$  converges to  $(x^*, \lambda^*)$  satisfies the second-order sufficient optimality condition. Suppose that the method uses a sequence of bounded matrices  $M_k$  such that  $H_k = (Z_k^{-T} M_k Z_k^{-1})^{-1}$  is positive definite and satisfies

$$(4.21) \quad \lim_{k \rightarrow \infty} \frac{||[Z_k^{-T} L(x^k, \lambda^{k+1}) Z_k^{-1} - H_k^{-1}] p^k||}{||x^{k+1} - x^k||} = 0$$

---

\* It is sufficient to assume that  $f$  and  $c$  are  $C^2$  differentiable and have Lipschitz continuous second derivatives.

with  $p^k = H_k g^k$ . Then the sequence  $\{x^k\}$  is TWO-STEPS SUPERLINEARLY converging to  $x^*$ , i.e.

$$(4.22) \quad \lim_{k \rightarrow \infty} \frac{\|x^{k+1} - x^*\|}{\|x^{k-1} - x^*\|} = 0. \quad \blacksquare$$

Proof : Notice first by combining (4.14) and (4.9) that

$$(4.23) \quad \frac{K_4}{K_1} \|x^k - x^*\| \leq \|x^{k+1} - x^k\| = \|d^k\| \leq \frac{K_3}{K_2} \|x^k - x^*\|;$$

hence

$$(4.24) \quad \frac{K_2}{K_3} \|d^{k+1}\| \leq \|x^{k+1} - x^*\| \leq \left(1 + \frac{K_3}{K_2}\right) \|x^k - x^*\|.$$

Combining (4.9) (for index  $k+1$ ), (4.15), (4.20) we obtain

$$(4.25) \quad \|x^{k+1} - x^*\| \leq K_1 \left[ \|[Z_k^{-T} L(x^k, \lambda^{k+1}) Z_k^{-1} - H_k^{-1}] p^k\| + (K_5 + K_6) \|d^k\|^2 + K_7 \|c^k\| \right];$$

using (4.23) and (4.24) (for index  $k-1$ ) yields

$$(4.26) \quad \frac{\|x^{k+1} - x^*\|}{\|x^{k-1} - x^*\|} \leq \frac{K_3}{K_2} \left(1 + \frac{K_3}{K_2}\right) \left\{ K_1 \frac{\|[Z_k^{-T} L(x^k, \lambda^{k+1}) Z_k^{-1} - H_k^{-1}] p^k\|}{\|x^{k+1} - x^k\|} + \right. \\ \left. + (K_5 + K_6) \|d^k\| \right\} + \frac{K_1 K_3 K_5 K_7}{K_2} \|d^{k-1}\|$$

Taking the limit as  $k \rightarrow \infty$  yields (4.22) given that  $\{d^k\} \rightarrow 0$  and (4.21) holds.  $\blacksquare$

Theorem 4.1 shows that the two-steps superlinear convergence of the general method (4.1) (4.2) depends only upon how the restriction of  $M_k$  to the subspace  $T_k$  approximates the similar restriction of the Hessian of the Lagrangian  $L(x^k, \lambda^{k+1})$ . Actually it only requires that this approximation be adequate along the direction  $p^k$ ; this result, mentioned by Powell for a particular method where the orthogonal restriction was envisaged, thus generalizes a similar condition given by Dennis-More (Ref.23) for Quasi-Newton methods for unconstrained minimization.

In other words it is not necessary in order to obtain two-step superlinear convergence to update the full  $n \times n$  matrix  $M_k$ . In our general framework, we are able to exploit fully this observation by considering matrices  $M_k$  satisfying (3.18) as discussed in Section 3.3. A simply choice consists in taking

$$(4.27) \quad M_k = Z_k^T G_k Z_k$$

where  $G_k$  is an  $(n-m) \times (n-m)$  positive definite symmetric matrix updated according to the modified generalized BFGS formula (3.23) ; notice that  $H_k = G_k^{-1}$  and can be updated directly according to (3.25). It remains to show that the update formula (3.23) verifies the Generalized Dennis-Moré condition (4.21) which requires very technical estimates ; we shall not attempt to prove this result in this paper (see Powell (Ref.9) where (4.21) is established for the updating of  $M_k$  by a formula analogous to (3.23)). We shall refer to such methods as Reduced Quasi-Newton Methods.

As discussed in Section 3.3, method (4.1) then uses the simplified Quasi-Newton direction

$$(4.28) \quad d^k = - Z_k^T H_k g^k - A_k^T c^k$$

and generates Lagrange multipliers

$$(4.29) \quad \lambda^{k+1} = - A_k^{-T} \nabla f(x^k) ;$$

notice that only  $H_k$  is required in the computations.

Finally, Theorem 4.1 can be interpreted in the framework of multipliers methods proposed in Section 3.4. It establishes the two-steps superlinear convergence of the diagonalized multipliers method where the minimization of the augmented Lagrangian

$$(4.30) \quad \min_y \ell(y, \mu^{k+1}, r)$$

is performed by one step of a Quasi-Newton method starting from  $x^k$  and using the approximate Hessian

$$(4.31) \quad N_k = Z_k^T G_k Z_k + r A_k^T A_k = S_k^T \begin{bmatrix} rI & 0 \\ 0 & G_k \end{bmatrix} S_k ,$$

and the multipliers are updated according to the formula

$$(4.32) \quad \begin{aligned} \mu^{k+1} &= - A_k^{-T} \nabla f(x^k) + r c^k \\ &= \mu^k + r c^k - A_k^{-T} \nabla_x \ell(x^k, \mu^k, r) . \end{aligned}$$

Such a result extends the analysis of Byrd (Ref.24) who considered diagonalized methods using Newton steps for (4.30).

#### 4.2. Superlinear Convergence of a Reduced Quasi-Newton Method with Feasibility Improvement.

Going back to the proof of Theorem 4.1 we can observe that superlinear convergence could be established had

$$\lim_{k \rightarrow \infty} \frac{\|c(x^k)\|}{\|x^k - x^*\|} = 0 .$$

To achieve this result we consider a modification of method (4.1) (4.2) where the next iterate is now given by

$$(4.33) \quad x^{k+1} = x^k + d^k + e^k$$

where  $d^k$  is still defined by (4.2) and  $e^k$  is an additional step to improve the error on the constraints, defined by

$$(4.34) \quad e^k = - A_k^{-T} c(x^k + d^k) ;$$

the additional step  $e^k$  is incorporated only if it actually improves feasibility, i.e. if given  $\alpha \in (0, \frac{1}{2})$ .

$$(4.35) \quad ||c(x^k + d^k + e^k)|| \leq (1-\alpha) ||c(x^k + d^k)||$$

By (4.34) we conclude using (4.15) that

$$(4.36) \quad ||e^k|| \leq K_8 ||d_k||^2 ,$$

which show that (4.35) is eventually satisfied as  $d^k \rightarrow 0$ .

We assume again that the sequence  $\{x^k\}$  defined by (4.33) and the associated sequence of Lagrange multipliers  $\{\lambda^k\}$  (still defined by (4.6)) converge to  $(x^*, \lambda^*)$  satisfying the second-order sufficient optimality conditions. We also assume that (4.4) and (4.5) still hold.

Propositions 4.1, 4.2, 4.3, 4.4 are obviously still valid together with estimates (4.9), (4.14), (4.15). It is easily verified, using (4.36) that the estimate (4.20) of  $||g^{k+1}||$  still hold, with possibly a larger constant  $K'_6$ .

To establish the rate of convergence we need an estimate of  $||c(x^{k+1})||$  for the new scheme. Taylor expansion of  $c(x^{k+1})$  leads now to

$$||c(x^{k+1})|| = ||c(x^k + d^k) + A_k e^k + \int_0^1 [c'(x^k + d^k + te^k) - c'(x^k)] e^k dt|| ;$$

hence, by the continuity of  $c''$ ,

$$(4.37) \quad ||c^{k+1}|| \leq K_9 ||d^k + e^k|| ||e^k|| .$$

THEOREM 4.2. Assume that  $f$  and  $c$  are  $C^\sigma$  differentiable with  $\sigma \geq 3$  and that the MODIFIED Quasi-Newton method (4.33) (4.2) (4.34) generates a sequence of approximate solutions  $\{x^k\}$  and a sequence of Lagrange multipliers  $\{\lambda^k\}$  (defined by (4.6)) converging to  $(x^*, \lambda^*)$  satisfying the second-order sufficient optimality condition.

Suppose that the method uses a sequence of bounded matrices  $M_k$  such that  $H_k = (Z_k^{-T} M_k Z_k)^{-1}$  is positive definite and satisfies

$$(4.38) \quad \lim_{k \rightarrow +\infty} \frac{||[Z_k^{-T} L(x^k, \lambda^{k+1}) Z_k^{-1} - H_k^{-1}] p^k||}{||x^{k+1} - x^k||} = 0 ,$$

with  $p^k = H_k s^k$ . Then the sequence  $\{x^k\}$  is SUPERLINEARLY converging to  $x^*$ , i.e.

$$(4.39) \quad \lim_{k \rightarrow +\infty} \frac{||x^{k+1} - x^*||}{||x^k - x^*||} = 0 . \quad \blacksquare$$

Proof : We still have

$$(4.40) \quad \frac{K_4}{K_1} ||x^k - x^*|| \leq ||d^k|| \leq \frac{K_3}{K_2} ||x^k - x^*|| ,$$

for the modified method, however,

$$(4.41) \quad ||x^{k+1} - x^k|| = ||d^k + e^k|| \leq ||d^k|| + K_8 ||d^k||^2 \leq K_{10} ||d^k|| ,$$

We thus obtain, using (4.25), (4.37), (4.40), (4.41)

$$(4.42) \quad \frac{||x^{k+1} - x^*||}{||x^k - x^*||} \leq \frac{K_3 K_{10}}{K_2} \left\{ K_1 \frac{||[Z_k^{-T} L(x^k, \lambda^{k+1}) Z_k^{-1} - H_k^{-1}] p^k||}{||x^{k+1} - x^k||} + \right. \\ \left. + K_7 K_9 ||e^{k-1}|| \right\} + \frac{K_3}{K_2} (K_5 + K_6') ||d^k||$$

Given (4.38) and the assumption that  $d^k \rightarrow 0$  (hence  $e^k \rightarrow 0$ ), we obtain (4.39), showing the superlinear convergence of the modified method.  $\blacksquare$

The comments following Theorem 4.1 are still in order and lead to the same strategy : we only need to update an  $(n-m) \times (n-m)$  positive definite matrix  $G_k = H_k^{-1}$  according to (3.23).

Notice that, starting from a feasible point  $\bar{x}^k$  (i.e.  $e^k = 0$ ), method (4.43) is equivalent to the Quasi-Newton method along geodesics of (Ref.1) where only one step of the restoration phase is performed (the stepsize being taken as 1).

Examples of choice of  $A_k^-$  and  $Z_k^-$  have been presented in (Ref.1) and include the partitioned right inverse formulae (3.15) and the Q.R. right inverse formulae (3.17).



## 5. STEPSIZE SELECTION

We now introduce a stepsize parameter  $t > 0$  and consider for the Reduced Quasi-Newton Method with Feasibility Improvement the parametrized arc of parabola starting from  $x^k$  and tangent to  $d^k$ , formally similar to the parabolic arc introduced in (Ref.1),

$$(5.1) \quad x(t) = x^k + t d^k + t^2 e^k ,$$

where  $d^k$  and  $e^k$  are respectively defined by (4.28) and (4.34) (if (4.35) is satisfied).

Following Han (Ref.10) we choose the stepsize  $t_k$ , defining the new iterate

$$(5.2) \quad x^{k+1} = x(t_k) ,$$

to achieve a sufficient decrease of the exact penalty function, analyzed in section 2,

$$(5.3) \quad \Phi(x, r_{k+1}) = f(x) + r_{k+1} \sum_{i=1}^m |c_i(x)| .$$

The non-decreasing sequence of penalty parameters is defined recursively by

$$(5.4) \quad r_{k+1} = \text{Max} \{ r_k, \text{Max}_i |\lambda_i^{k+1}| \} ,$$

starting from  $r_0 > 0$  ; the  $m$ -dimensional vector  $\lambda^{k+1}$  is taken as the Lagrange multipliers vector at the solution of the quadratic programming problem (3.4) and given by (4.29).

Instead of requiring  $t_k$  to achieve an approximate minimization of the form proposed in (Ref.10), we follow the spirit of Powell (Ref.20) and select  $t_k = 2^{-\ell}$  for the first index  $\ell$  of the sequence  $\{0, 1, 2, \dots\}$  such that

$$(5.5) \quad \Phi(x(t_k), r_{k+1}) \leq \Phi(x^k, r_{k+1}) - \alpha t_k \Psi(x^k, d^k, r_{k+1})$$

with  $\alpha \in (0, \frac{1}{2})$  and  $\Psi(x, d, r)$  defined by

$$(5.6) \quad \Psi(x, d, r) = r \sum_{i=1}^m |c_i(x)| - \langle \nabla f(x), d \rangle.$$

PROPOSITION 5.1. Given  $x^k$ , let  $d^k$ ,  $\lambda^{k+1}$  and  $r_{k+1}$  be defined by (4.28), (4.29) and (5.4) ; then

$$(5.7) \quad \Psi(x^k, d^k, r_{k+1}) \geq 0,$$

where equality holds iff  $(x^k, \lambda^{k+1})$  satisfies the first-order optimality conditions (2.4). ■

Proof : Notice that, using (4.28) and (4.29), formula (5.6) yields

$$(5.8) \quad \Psi(x^k, d^k, r_{k+1}) = \langle g^k, H_k g^k \rangle + r_{k+1} \sum_{i=1}^m |c_i(x^k)| - \langle \lambda^{k+1}, c(x^k) \rangle ;$$

inequality (5.7) results from the choice (5.4) of  $r_{k+1}$  and the positive definiteness of  $H_k$ . ■

The stepsize selection rule (5.5) thus insures a sufficient decrease of the exact penalty function from a non-critical iterate  $x^k$ . The convergence analysis of the algorithm must however distinguish between two situations.

THEOREM 5.1. (Global Convergence) Assume that  $f$  and  $c$  are  $C^\sigma$  differentiable with  $\sigma \geq 2$  and that  $c$  is a submersion. If the sequence  $\{r_k\}$  defined by (5.4) increases infinitely, then the sequence  $\{x^k\}$ , constructed by (5.1), (5.2) (5.5), has no accumulation point ; if  $r_k$  is increased only a finite number of times according to (5.4), then any accumulation point of the sequence  $\{x^k, \lambda^{k+1}\}$  satisfies the first-order optimality conditions. ■

Proof : a) Suppose that  $r_k \rightarrow +\infty$  as  $k \rightarrow +\infty$  and that the sequence  $\{x^k\}$  has an accumulation point  $x^*$ , i.e. there exists a subsequence  $x_{k_i} \rightarrow x^*$  as  $i \rightarrow +\infty$ .

Define

$$(5.9) \quad r(x) = \max_{1 \leq i \leq m} |\lambda_i(x)| ,$$

where  $\lambda(x)$  is the Lagrange multipliers vector, defined by

$$(5.10) \quad \lambda(x) = -A_x^{-T} \nabla f(x) .$$

Since  $r$  is a continuous function, given  $\epsilon > 0$ ,

$$(5.11) \quad r^* = \max_{x \in \bar{B}(x^*, \epsilon)} r(x) < +\infty .$$

There also exists  $N > 0$  such that  $x_{k_i} \in B(x^*, \epsilon)$  for all  $i \geq N$ .

Since  $r_k$  is increased infinitely often according to (5.4), we must have

$$(5.12) \quad r_{k_i} < r(x_{k_i}) \leq r^* < +\infty \quad \forall i \geq N ,$$

which contradicts the fact that  $r_k \rightarrow +\infty$ .

b) Suppose now that  $r_k$  is increased finitely many times, i.e. that there exist  $r > 0$  and an integer  $N$  such that

$$(5.13) \quad r_{k+1} = r \quad \forall k \geq N ,$$

which implies that the Lagrange multipliers remain bounded :

$$(5.14) \quad |\lambda_i^{k+1}| \leq r \quad \text{for } i = 1, \dots, m , \quad \forall k \geq N .$$

Let  $(x^*, \lambda^*)$  an accumulation point of the sequence  $\{x^k, \lambda^{k+1}\}$ ; we assume for simplicity that  $x^k \rightarrow x^*$ . Suppose that  $(x^*, \lambda^*)$  does not satisfy the first-order optimality conditions (2.4); then by Proposition 5.1

$$(5.15) \quad \Psi(x^*, d^*, r) = \delta > 0 .$$

There exists  $N' \geq N$  such that

$$(5.16) \quad \Psi(x^k, d^k, r) > \frac{\delta}{2} \quad \forall k \geq N'.$$

We can evaluate  $\Phi(x(t), r)$  along the parabolic arc (5.1) starting from  $x^k$ ,  $k \geq N'$ , for  $t \in [0, 1]$  :

$$(5.17) \quad \Phi(x(t), r) = f(x^k + t d^k + t^2 e^k) + r \sum_{i=1}^m |c_i(x^k + t d^k + t^2 e^k)|$$

Using the definitions (4.28) and (4.34) of  $d^k$  and  $e^k$ , second-order Taylor expansions of  $f$  and  $c_i$  yield the majorization

$$(5.18) \quad \Phi(x(t), r) \leq \Phi(x^k, r) - t \Psi(x^k, d^k, r) + t^2 \theta(x^k, d^k, e^k, r),$$

where  $\theta(x^k, d^k, e^k, r)$  is a positive bounded term since  $\nabla^2 f$  and  $\nabla^2 c_i$  are continuous hence bounded on  $\bar{B}(x^*, \varepsilon)$  :

$$(5.19) \quad 0 < \theta(x^k, d^k, e^k, r) \leq M.$$

There exists an integer  $L \geq 0$  such that

$$2^{-L} \leq \frac{(1-\alpha)\delta}{2M} < 2^{-L+1}$$

provided  $0 < \alpha < 1$  and  $\delta \leq 4M$  ; then  $t = 2^{-L}$  is an admissible stepsize since

$$(5.20) \quad \Phi(x(2^{-L}), r) \leq \Phi(x^k, r) - \alpha 2^{-L} \Psi(x^k, d^k, r).$$

Thus the selection rule (5.5) defines a sequence of stepsizes  $t_k$  bounded from below

$$t_k \geq 2^{-L} \quad \forall k \geq N'.$$

By definition (5.2) of the new iterate  $x^{k+1}$ , we have

$$(5.21) \quad \Phi(x^{k+1}, r) - \Phi(x^k, r) \leq -\alpha t_k \Psi(x^k, d^k, r) \leq -\alpha \delta 2^{-L-1},$$

which implies that  $\Phi(x^k, r) \rightarrow -\infty$ , in contradiction with the continuity of  $\Phi(., r)$  which guarantees that  $\Phi(x^k, r) \rightarrow \Phi(x^*, r)$ .

Hence  $(x^*, \lambda^*)$  must satisfy the first-order optimality conditions. ■

We assume now that  $f$  and  $c_i$  are  $C^\sigma$  differentiable functions with  $\sigma \geq 3$ . The superlinear convergence established in Section 4 holds since we can establish that the stepsize  $t_k = 1$  satisfies (5.5) after a finite number of iterations.

PROPOSITION 5.2. There exists an integer  $N > 0$  such that the stepsize  $t_k = 1$  satisfies the selection rule (5.5) for  $k \geq N$ . ■

Proof : Third-order Taylor expansion of  $f$  yields, using (4.36),

$$(5.22) \quad f(x^k + d^k + e^k) = f(x^k) + \langle \nabla f(x^k), d^k \rangle + \langle \nabla f(x^k), e^k \rangle \\ + \frac{1}{2} \langle d^k, \nabla^2 f(x^k) d^k \rangle + \mathcal{O}(\|d^k\|^3).$$

Combining (4.36) and (4.37), we obtain

$$(5.23) \quad \|c(x^k + d^k + e^k)\| = \mathcal{O}(\|d^k\|^3),$$

hence

$$(5.24) \quad \Phi(x^k + d^k + e^k, r) = \Phi(x^k, r) - \Psi(x^k, d^k, r) + \langle \lambda^{k+1}, c(x^k + d^k) \rangle \\ + \frac{1}{2} \langle d^k, \nabla^2 f(x^k) d^k \rangle + \mathcal{O}(\|d^k\|^3),$$

where we have made use of the definitions (4.34) and (4.29) of  $e^k$  and  $\lambda^{k+1}$ .

Using the third-order Taylor expansions of  $c_i$ , we obtain

$$(5.25) \quad \Phi(x^k + d^k + e^k, r) = \Phi(x^k, r) - \Psi(x^k, d^k, r) + \frac{1}{2} \langle d^k, L(x^k, \lambda^{k+1}) d^k \rangle + \mathcal{O}(\|d^k\|^3)$$

With the expressions (5.8) of  $\Psi(x^k, d^k, r)$  and (4.28) of  $d^k$ , we can rewrite (5.25) as

$$(5.26) \quad \Phi(x^k + d^k + e^k, r) - \Phi(x^k, r) + \alpha \Psi(x^k, d^k, r) \leq -(\frac{1}{2} - \alpha) \Psi(x^k, d^k, r) \\ + \frac{1}{2} \langle p^k, [Z_k^{-T} L(x^k, \lambda^{k+1}) Z_k^{-1} - H_k^{-1}] p^k \rangle + O(\|d^k\|^3)$$

with  $p^k = H_k g^k$ . Notice that if the sequence  $H_k$  generated by the algorithm satisfies (4.38) the second term of the RMS is  $O(\|d^k\|^2)$  while, by (5.8),

$$(5.27) \quad \Psi(x^k, d^k, r) = O(\|d^k\|^2) ;$$

hence, provided  $\alpha < \frac{1}{2}$ , (5.26) shows that

$$(5.28) \quad \Phi(x^k + d^k + e^k, r) \leq \Phi(x^k, r) - \alpha \Psi(x^k, d^k, r)$$

once  $\|d^k\|$  is sufficient small (i.e. after a finite number of iterations since  $\{x^k\} \rightarrow x^*$  satisfying the first-order optimality conditions). ■

We have thus established the global and superlinear convergence of the Reduced Quasi-Newton Method with Feasibility Improvement which we summarize as follows :

Given  $x^k \in \mathbb{R}^n$ ,  $H_k$  positive definite,  $A_k^-$ ,  $Z_k^-$  satisfying (3.6),  $\alpha \in (0, \frac{1}{2})$ ,  $r_k > 0$  ;

- i) compute the constraints residues :  $c^k = c(x^k)$  ;
- ii) compute the reduced gradient :  $g^k = Z_k^{-T} \nabla f(x^k)$  ;
- iii) compute the Quasi-Newton direction  $d^k = -Z_k^- H_k g^k - A_k^- c^k$  ;
- iv) compute the feasibility improvement direction :

$$\text{let } \tilde{x}^k = x^k + d^k \text{ and } e^k = -A_k^- c(\tilde{x}^k) ;$$

$$\text{if } \|c(\tilde{x}^k + e^k)\| > (1-\alpha) \|c(\tilde{x}^k)\| \text{ then } e^k \leftarrow 0 ;$$

- v) penalty parameter :

$$\text{compute the Lagrange multipliers } \lambda^{k+1} = -A_k^{-T} \nabla f(x^k)$$

$$\text{let } r_{k+1} = \text{Max} \{r_k, \text{Max}_i |\lambda_i^{k+1}|\}$$

vi) stepsize selection : let  $\ell$  be the smallest integer such that

$$\Phi(x^k + 2^{-\ell} d^k + 2^{-2\ell} e^k, r_{k+1}) \leq \Phi(x^k, r_{k+1}) - \alpha 2^{-\ell} \Psi(x^k, d^k, r_{k+1}) ;$$

$$\text{let } x^{k+1} = x^k + 2^{-\ell} d^k + 2^{-2\ell} e^k .$$

## 6. CONCLUSIONS

In this paper we have presented a superlinearly and globally convergent algorithm for the minimization of a differentiable function over a differential manifold. This method can be viewed as a particularly efficient approximation of the Quasi-Newton method along geodesics presented in (Réf. 1) where the feasibility of the successive iterates is not enforced. It is also related to the class of diagonalized multipliers methods (Ref. 19) and to the increasingly popular variable metric methods for constrained optimization (Ref. 8, 9) but offers several advantages : from a computation viewpoint it only requires to update a reduced-size approximate Hessian of the Lagrangian (at the Montreal meeting, both Abadie and Tanabe presented algorithms with a similar feature) ; on the theoretical side the method is guaranteed to converge with a superlinear rate (a similar procedure was presented by Mayne and Polak (Ref. 25) at the Montreal meeting). Additional analysis is needed to extend our method to mathematical programming problems with inequality constraints (see Ref. 1 §.6 for a possible approach).

# REFERENCES

1. D. GABAY, "Minimizing a Differentiable Functional over a Differential Manifold. Part I : Descent Methods along Geodesics and Practical Implementation", paper presented at the 10th International Symposium on Mathematical Programming, Montreal (1979).
2. H. MUKAI and E. POLAK, "On the Use of Approximations in Algorithms for Optimization Problems with Equality and Inequality Constraints", SIAM J. Numer. Anal. 15, pp. 674-693 (1978).
3. R.T. ROCKAFELLAR, "Augmented Lagrange Multiplier Functions and Duality in Nonconvex Programming", SIAM J. Control 12, pp. 555-562 (1973).
4. D.P. BERTSEKAS, "Multiplier Methods : a survey", Automatica 12, pp. 133-145 (1976).
5. M.C. BIGGS, "Constrained Minimization using Recursive Equality Quadratic Programming", in Numerical methods for Nonlinear Optimization, F.A. Lootsma, Ed., Academic Press, London, pp. 411-428 (1972).
6. S.M. ROBINSON, "A Quadratically Convergent Algorithm for General Nonlinear Programming Problems", Mathematical Programming 3, 145-156 (1972).
7. U.M. GARCIA-PALOMARES and O.L. MANGASARIAN, "Superlinearly Convergent Quasi-Newton Algorithms for Nonlinearly Constrained Optimization Problems", Mathematical Programming 11, pp. 1-13 (1976).
8. S.P. HAN, "Superlinearly Convergent Variable Metric Algorithms for General Nonlinear Programming Problems", Mathematical Programming 11, pp. 263-282 (1976).



9. M.J.D. POWELL, "The convergence of Variable Metric Methods for Nonlinearly Constrained Optimization Calculations", in Nonlinear Programming 3, O.L. Mangasarian, R.R. Meyer and S.M. Robinson eds., Academic Press, New-York, pp. 27-63 (1978).
10. S.P. HAN, "A Globally Convergent Method for Nonlinear Programming", J. Opt. Theory Appl. 22, pp. 297-309 (1977).
11. A. MIELE, J. TIETZE, A.V. LEVY, "Summary and Comparison of Gradient Restoration Algorithms for Optimal Control, J. Opt. Theory Appl. 10, pp. 381-403 (1972).
12. K. TANABE, "Differential Geometric Methods in Nonlinear Programming", in Applied Nonlinear Analysis, V. Lakshmikanthan ed., Academic Press, New-York, pp. 707-719 (1979).
13. J.W. MILNOR, Topology From the Differentiable Viewpoint, University Press of Virginia, Charlottesville (1965).
14. M.W. HIRSCH, Differential Topology, Springer, Heidelberg and New-York (1976).
15. D.G. LUENBERGER, Introduction to Linear and Nonlinear Programming, Addison-Wesley, Reading, Mass. (1973).
16. R.B. WILSON, "A Simplicial Method for Convex Programming", Ph.D dissertation, Harvard University, Cambridge, Mass. (1963).
17. B. NOBLE, Applied Linear Algebra, Prentice Hall, Englewoods-Cliffs, N.J. (1969).
18. M.J.D. POWELL, "Algorithms for Nonlinear Constraints that use Lagrangian Functions", Mathematical Programming 14, pp. 224-248 (1978).

19. R.A. TAPIA, "Diagonalized Multiplier Methods and Quasi-Newton Methods for Constrained Optimization", J. Opt. Theory Appl. 22, pp. 135-194 (1977).
20. M.J.D. POWELL, "A Fast Algorithm for Nonlinearly Constrained Optimization Calculations", in Numerical Analysis, G.A. Watson ed., Springer, Heidelberg and New-York, pp. 144-157 (1978).
21. A. MIELE, E.E. CRAGG and A.V. LEVY, "Use of the Augmented Penalty Function in Mathematical Programming Problems, Part. II, J. Opt. Theory Appl. 8, pp. 336-349 (1971).
22. R.A. TAPIA, "Quasi-Newton Methods for Equality Constrained Optimization : Equivalence of Existing Methods and a New Implementation", paper presented at Nonlinear Programming Symposium 3, Madison, Wis. (1977).
23. J.E. DENNIS and J.J. MORE, "Quasi-Newton Methods, Motivation and Theory", SIAM Review 19, pp. 46-89 (1977).
24. R.H. BYRD, "Local Convergence of the Diagonalized Method of Multipliers", J. Opt. Theory Appl. 26, pp. 485-500 (1978).
25. D.Q. MAYNE and E. POLAK, "A Superlinearly Convergent Algorithm for Constrained Optimization Problems", paper presented at 10th International Symposium on Mathematical Programming, Montreal (1979).

